

Tartalomjegyzék

1. Bevezetés	3
2. Broyden életrajza	4
3. Mátrixok alapvető jellemzése	5
3.1. Vektorok és mátrixok	5
3.2. Lineáris függetlenség	10
3.3. Nemszinguláris mátrixok	12
3.4. Vektornormák és mátrixnormák	15
4. Lineáris egyenletrendszerek megoldása, osztályozása	17
4.1. Direkt eljárások	18
4.2. Iteratív eljárások	23
4.3. Sajátértékek és sajátvektorok	24
4.4. Kovergens mátrixok	32
5. Felső relaxációs eljárás (SOR)	34
5.1. Általános konvergencia eredmények	34
5.2. Szimmetrikus mátrixok konvergenciája	37
5.3. Nemszimmetrikus mátrixok konvergenciája	37
5.4. Konklúzió	39
6. A konjugált gradiens módszerek új rendszertana	40
6.1. A témakör felvezetése	40
6.2. Szimmetrikus mátrixok	45
6.3. Nemszimmetrikus mátrixok	46
6.4. Konklúzió	50
7. Blokk konjugált gradiens módszer (BICG)	51
7.1. A Lanczos-féle és a Hestenes-Stiefel-féle algoritmusok	51
7.2. Az összeomlás elkerülésének feltételei	52
7.3. Konklúzió	54
8. Matlab tesztfeladatokon való összehasonlítás	55
8.1. Teszt-specifikáció	55
8.2. A tesztfeladatok elemzése	56
8.3. Teszteredmények	57
8.4. Konklúzió	59
9. Összefoglaló	60

10.Summary	61
11.Függelék	63
11.1. A fő program	63
Hivatkozások	65

1. Bevezetés

Dr. Abaffy József tanár úr három témát hirdetett meg Charles George Broyden barátja emlékére a 2017/18/1 félévben. Ezek a következők voltak:

1. Broyden és a feltétel nélküli optimalizálás,
2. Broyden és a nemlineáris egyenletrendszerek megoldása,
3. Broyden és a lineáris egyenletrendszerek megoldása.

Én a harmadik témát választottam. A lineáris algebra szerepe kiemelkedően fontos rengeteg területen, például jelentős geometriai, fizikai és mérnöki alkalmazásokkal rendelkezik, de a modern társadalomtudományokban is alkalmazzák. Számos más területen is találkozhatunk lineáris algebrával, szinte minden tudományág tartalmaz olyan modelleket, amelyek lineáris egyenletrendszerek megoldására vezetnek vissza valamilyen probléma megoldását.

Broyden-en kívül természetesen más szerzők műveiből is merítünk (lásd a hivatkozásoknál), de mivel a témakör igen nagy, teljességre nem törekedhettünk e téren. A diplomamunka pl. a legtöbb tétel bizonyítását tartalmazza, de van, ahol csak hivatkozunk a bizonyításokra. Vannak olyan eredmények, amiket a hely és az idő végeessége miatt teljesen kihagytunk, de természetesen ezek a kevésbé fontos eredmények.

A diplomamunka Broyden rövid életrajzával kezdődik. A 3. és a 4. fejezet a további fejezeteket alapozzák meg. Ezekben a fejezetekben a további fejezetekhez szükséges fogalmakat vezetünk be. Definíciókat adunk meg és tételeket mondunk ki. Tisztázzuk a jelöléseket. A témakör felvezetése annyiban különleges, hogy egyáltalán nem használjuk a determináns fogalmát. Ebben Broyden-hez igazodunk, mivel ezek a fejezetek nagyban támaszkodnak a Basic Matrices [3] című könyvére. A fejezetek struktúrája szintén kicsit eltér a megszokottól. Ebben szintén Broyden-hez igazodtunk. A témakör más irányból való megközelítése számomra igen hasznos volt. Az 5-7. fejezetek Broyden, a numerikus lineáris algebra témakörében elért eredményeit foglalják össze. Az 5. fejezetben vizsgáljuk a felső relaxációs eljárás konvergencia kritériumait. A 6. fejezetben egy átfogó rendszertant adunk a konjugált gradiens módszerekre. A 7. fejezetben a blokk konjugált gradiens módszert vizsgáljuk. A 8. fejezetben MATLAB tesztfeladatokon hasonlítjuk össze a pcg, SYMMLQ, bicgstabl, QMR és az LSQR módszereket.

Szeretném megköszönni Dr. Abaffy József segítségét a diplomamunka írásával kapcsolatban.

2. Broyden életrajza



Charles George Broyden (1933. február 3. - 2011. május 20.) szerény családi háttérrel, Angliában született. Édesapja gyári munkásként, édesanyja háztartásbeliként dolgozott. Szülei ennek ellenére a kezdetektől tanulásra biztatták. Charles már gyerekként rengeteget olvasott, az iskolában jól teljesített. Édesapja sajnálatos módon meghalt tuberkulózisban, mikor Charles még csak 11 éves volt. Ez még jobban megnehezítette családi helyzetüket, de édesanyja így is arra biztatta, hogy egyetemre menjen. A King's College London egyetemen szerzett fizikus diplomát 1955-ben. A következő 10 évet az iparban töltötte. Ezután 1965-1967 között az Aberystwyth Egyetemen tanított, majd a University of Essex egyetemen 1967-ben professzor, később a matematika intézet dékánja lett. 1986-ban innen visszavonult, 1990-ben a Bolognai Egyetemen fogadott el professzori kinevezést. Jelentős szerepe volt a kvázi-Newton módszerek kifejlesztésében. A kvázi-Newton módszerek előtt a nemlineáris optimalizálási problémákat gradiens alapú módszerekkel oldották meg. Nemlineáris esetben az ehhez szükséges Hessian mátrix kiszámítása legtöbbször nem praktikus. A kvázi-Newton módszerek kifejlesztésére irányuló munka az 1960-as és 1970-es években zajlott, a nemlineáris optimalizálás egy izgalmas időszakában. A kutatásban részt vett még például Bill Davidon, Roger Fletcher és Mike Powell. A kutatások eredménye valódi ipari alkalmazások problémáinak megoldására adott eszközöket.

Az iparban töltött éveit alatt Broyden a Davidon-Fletcher-Powell (DFP) módszert adaptálta nemlineáris problémákra. Ez vezetett az 1965-ös klasszikus "A class of methods for solving nonlinear simultaneous equations" cikkéhez a Mathematics of Computation folyóiratban. Ez a munka széleskörűen elismert, mint a 20. század egyik legnagyobb numerikus analízis eredménye. A University of Essex egyetemen a DFP módszer optimalizálásra fókuszált. Feltűnt neki, hogy habár a módszer jól működik, néha furcsa eredményeket produkál. Kerekítési hibákra gyanakodott. A kutatása 1970-ben egy új, továbbfejlesztett módszerhez vezetett. Tőle függetlenül, nagyjából egy időben, Fletcher, Donald Goldfarb, és David Shanno is ugyanerre az eredményre jutott. Ezért az új módszert a neveik kezdőbetűiből BFGS módszernek

nevezték el. Más kutatások folytatták a kvázi-Newton módszerek optimalizálását, de a BFGS módszer még ma is a leginkább választott, ha a Hessian mátrix kiszámítása túl költséges. 1981-től Abaffy Józseffel és Emilio Spedicato-val az ABS (Abaffy-Broyden-Spedicato) módszereken dolgozott.

Később a numerikus lineáris algebrára fókuszált, ezen belül is a konjugált gradiens módszerekre és ezek rendszertanára. A diplomamunka a kutatásainak a konjugált gradiens módszerekkel kapcsolatos részét dolgozza fel. 2011-ben 78 évesen, egy szívroham komplikációiba belehalt.

Feleségével, Joan-nal, 1959-ben házasodtak össze. Négy gyerekük született, a legidősebb, Robbie, sajnos 4 éves korában meghalt. Broyden nagy örömét lelte családjában, gyermekeiben, Christopher-ben, Jane-ben és Nicholas-ban. Szeretett mádárlesre járni, zenével foglalkozni, kórusban énekelni, vitorlázni. A helyi közösség és az egyházi közösség aktív tagja volt. Hét unokája született. A legidősebb unokája, Tom, az Oxford egyetemen tanult matematikát, a második legidősebb unokája, Matt, a Warwick Egyetemen tanul matematikát, Ben pedig mérnöknek tanul a Swansea Egyetemen. Így nagyapjuk nyomában járnak.

Köszönet Joan Broyden-nek a életrajzban nyújtott segítségéért.

3. Mátrixok alapvető jellemzése

Ebben a fejezetben a további fejezetekhez szükséges fogalmakat vezetünk be. Definíciókat adunk meg és tételeket mondunk ki. Tisztázzuk a jelöléseket. [3] [13]

3.1. Vektorok és mátrixok

Definíció. A valós vektor a valós számok egy rendezett halmaza.

Definíció. Egy vektor elemeinek a száma a vektor rendje, vagy más szóval a vektor dimenziója.

A vektorokat a diplomamunkában vastag kisbetűvel jelöljük. Például $\mathbf{x} = [x_i]$, ahol x_i a vektor i -edik elemét jelöli. Oszlopvektoron vektort, sorvektoron vektor transzponáltat értünk. Az \mathbf{x} vektor transzponáltját \mathbf{x}^T -vel jelöljük.

Példa. Ha $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$, akkor $\mathbf{x}^T = \begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix}$.

Definíció. Legyenek $\mathbf{x} = [x_i]$, $\mathbf{y} = [y_i]$ és $\mathbf{z} = [z_i]$ n -ed rendű vektorok. Legyen $\mathbf{z} = \mathbf{x} + \mathbf{y}$. Ekkor $z_i = x_i + y_i$.

Definíció. Legyen $\mathbf{x} = [x_i]$ és $\mathbf{y} = [y_i]$ n -ed rendű valós vektor. A belső szorzata, vagy más néven skaláris szorzata a sorvektor \mathbf{x}^T -nek és az oszlopvektor \mathbf{y} -nak

$$\mathbf{x}^T \mathbf{y} = \sum_{i=1}^n x_i y_i. \quad (3.1)$$

Definíció. Az \mathbf{x} és \mathbf{y} vektorok egymásra ortogonálisak, ha a belső szorzatuk 0.

Definíció. Legyenek \mathbf{p}_i vektorok és y_i skalárok, $i = 1, 2, \dots, n$. A

$$\sum_{i=1}^n \mathbf{p}_i y_i \quad (3.2)$$

vektor a \mathbf{p}_i vektorok lineáris kombinációja.

Definíció. A valós mátrix azonos rendű valós vektorok egy rendezett halmaza.

A mátrixokat a diplomamunkában vastag nagybetűvel jelöljük. Egy $m \times n$ -es mátrixot értelmezhetünk úgy, mint m darab sorvektor, vagy mint n darab oszlopvektor. Röviden azt mondjuk, hogy a mátrix sorai és oszlopai. Az $\mathbf{A} = [a_{ij}]$ jelölésben a_{ij} az \mathbf{A} mátrix i -edik sorának j -edik eleme. A csupa nullából álló mátrixot vagy vektort $\mathbf{0}$ -val jelöljük. Egy mátrixra akkor mondjuk, hogy ritka, ha viszonylag kevés nullától különböző eleme van. Egy mátrixra akkor mondjuk, hogy kitöltött, ha viszonylag kevés nulla eleme van.

Példa. Legyen $\mathbf{A} = [a_{ij}]$ $m \times n$ -es mátrix. Ekkor $i = 1, 2, \dots, m$ és $j = 1, 2, \dots, n$.

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix}$$

Definíció. Ha egy mátrixnak ugyanannyi sora és oszlopa van, négyzetes mátrixnak nevezzük.

Definíció. A négyzetes mátrix rendje az oszlopainak száma.

Definíció. Az $\mathbf{A} = [a_{ij}]$ mátrix transzponáltját \mathbf{A}^T -vel jelöljük, jelentése $\mathbf{A}^T = [a_{ji}]$, azaz a mátrix sorainak és oszlopainak felcserélésével kapott mátrix.

Definíció. Legyen $\mathbf{A} = [a_{ij}]$ mátrix. A mátrix diagonális elemei azok az elemek, ahol $i = j$.

Definíció. Az $n \times n$ -es $\mathbf{D} = [d_{ij}]$ mátrix diagonális mátrix, ha $d_{ij} = 0$ minden $i \neq j$ indexre. A diagonális mátrixot $\text{diag}(d_i)$ -vel jelöljük, ahol d_i az i -edik diagonális elem.

Definíció. A $\text{diag}(d_i)$ mátrixot, ahol $d_i = 1$ minden i -re, egységmátrixnak nevezzük, és \mathbf{I} -vel jelöljük.

Definíció. Legyenek $\mathbf{A} = [a_{ij}]$, $\mathbf{B} = [b_{ij}]$ és $\mathbf{C} = [c_{ij}]$ $m \times n$ -es mátrixok. Legyen $\mathbf{C} = \mathbf{A} + \mathbf{B}$. Ekkor $c_{ij} = a_{ij} + b_{ij}$.

Definíció. Legyen \mathbf{A} $m \times n$ -es mátrix. Legyen $\mathbf{x} = [x_i]$ n -ed rendű vektor, $\mathbf{y} = [y_i]$ m -ed rendű vektor. Jelölje \mathbf{a}_i^T az \mathbf{A} mátrix i -edik sorát. Az \mathbf{A} mátrix és \mathbf{x} vektor mátrix-vektor szorzata $\mathbf{y} = \mathbf{Ax}$ és

$$y_i = \mathbf{a}_i^T \mathbf{x}. \quad (3.3)$$

Definíció. Legyen \mathbf{A} $m \times n$ -es mátrix és jelölje \mathbf{a}_i^T az \mathbf{A} mátrix i -edik sorát. Legyen \mathbf{B} $n \times p$ -s mátrix és jelölje \mathbf{b}_j a \mathbf{B} mátrix j -edik oszlopát. Az \mathbf{AB} mátrixszorzat eredménye a $\mathbf{C} = [c_{ij}]$ $m \times p$ -s mátrix, ahol

$$c_{ij} = \mathbf{a}_i^T \mathbf{b}_j. \quad (3.4)$$

Az \mathbf{AB} mátrixszorzat esetén az \mathbf{A} mátrix oszlopainak száma meg kell egyezzen a \mathbf{B} mátrix sorainak számával, hogy a megfelelő sorvektorok és oszlopvektorok belső szorzatai jól definiáltak legyenek. A mátrixszorzás általában nem kommutatív, $\mathbf{AB} \neq \mathbf{BA}$. Egy mátrixot az egységmátrixszal szorozva önmagát kapjuk, $\mathbf{AI} = \mathbf{IA} = \mathbf{A}$.

Definíció. Az \mathbf{A} mátrix idempotens, ha $\mathbf{A}^2 = \mathbf{A}$.

Definíció. Az \mathbf{A} mátrix szimmetrikus, ha $\mathbf{A} = \mathbf{A}^T$. Az \mathbf{A} mátrix antiszimmetrikus, ha $\mathbf{A} = -\mathbf{A}^T$.

Mivel $(\mathbf{A}^T)^T = \mathbf{A}$ és $(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T$, ezért $(\mathbf{A}^T \mathbf{A})^T = \mathbf{A}^T (\mathbf{A}^T)^T = \mathbf{A}^T \mathbf{A}$. Azaz az $\mathbf{A}^T \mathbf{A}$ valós mátrix mindig szimmetrikus.

Definíció. Legyen $\mathbf{x}^T = [\mathbf{x}_1^T \ \mathbf{x}_2^T]$ n -ed rendű vektor, ahol $\mathbf{x}_1^T = [x_1, x_2, \dots, x_r]$ és $\mathbf{x}_2^T = [x_{r+1}, x_{r+2}, \dots, x_n]$ és $1 \leq r < n$. Az \mathbf{x}_1 és \mathbf{x}_2 vektorok az \mathbf{x} vektor részvektorai, vagy partíciói.

Példa. Legyenek $\mathbf{x} = [\mathbf{x}_1 \ \mathbf{x}_2]$, $\mathbf{y} = [\mathbf{y}_1 \ \mathbf{y}_2]$, \mathbf{u} n -ed rendű vektorok. Legyen $\mathbf{u} = \mathbf{x} + \mathbf{y}$. Ha $\mathbf{x}_1 = [x_1, x_2, \dots, x_r]$, $\mathbf{x}_2 = [x_{r+1}, x_{r+2}, \dots, x_n]$, $\mathbf{y}_1 = [y_1, y_2, \dots, y_r]$, $\mathbf{y}_2 = [y_{r+1}, y_{r+2}, \dots, y_n]$, $1 \leq r < n$, akkor $\mathbf{u} = [\mathbf{x}_1 + \mathbf{x}_2 \ \mathbf{y}_1 + \mathbf{y}_2]$. Igaz az is, hogy $\mathbf{x}^T \mathbf{y} = \mathbf{x}_1^T \mathbf{y}_1 + \mathbf{x}_2^T \mathbf{y}_2$.

Mátrixokat is részmátrixokra particionálhatunk. Ennek nagy jelentősége, hogy nagy mátrixokat egyszerűbben kezelhetünk.

Definíció. Legyen \mathbf{A} $m \times n$ -es mátrix. Az \mathbf{A} mátrixból annak k számú sora ($1 \leq k \leq m-1$) és l számú oszlopa ($1 \leq l \leq n-1$) törlésével előállított $(m-k) \times (n-l)$ -es részmátrixot az \mathbf{A} mátrix minormátrixának nevezzük.

A művelettartás a minormátrixokra két particionált mátrix között nem feltétlen jól definiált.

Definíció. Particionált mátrixok egy halmaza egy művelet szerint jól particionált, ha a mátrixok minormátrixai között a művelet jól definiált.

Példa. Az \mathbf{E}_1 és \mathbf{F}_1 mátrixok az összeadás szerint jól particionáltak.

$$\mathbf{E}_1 = \left[\begin{array}{cc} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{array} \right] = \frac{\begin{array}{ccc|c} a_{11} & a_{12} & a_{13} & b_{11} \\ a_{21} & a_{22} & a_{23} & b_{21} \\ c_{11} & c_{12} & c_{13} & d_{11} \\ c_{21} & c_{22} & c_{23} & d_{21} \end{array}}{\begin{array}{ccc|c} j_{11} & j_{12} & j_{13} & k_{11} \\ j_{21} & j_{22} & j_{23} & k_{21} \\ l_{11} & l_{12} & l_{13} & m_{11} \\ l_{21} & l_{22} & l_{23} & m_{21} \end{array}}, \quad \mathbf{F}_1 = \left[\begin{array}{cc} \mathbf{J} & \mathbf{K} \\ \mathbf{L} & \mathbf{M} \end{array} \right] = \frac{\begin{array}{ccc|c} j_{11} & j_{12} & j_{13} & k_{11} \\ j_{21} & j_{22} & j_{23} & k_{21} \\ l_{11} & l_{12} & l_{13} & m_{11} \\ l_{21} & l_{22} & l_{23} & m_{21} \end{array}}{\begin{array}{ccc|c} j_{11} & j_{12} & j_{13} & k_{11} \\ j_{21} & j_{22} & j_{23} & k_{21} \\ l_{11} & l_{12} & l_{13} & m_{11} \\ l_{21} & l_{22} & l_{23} & m_{21} \end{array}}$$

A szorzás szerint \mathbf{E}_1 és \mathbf{F}_1 nem jól particionáltak, mert a megfelelő minormátrixok szorzata (3.4) szerint nem képezhető.

Példa. Az \mathbf{E}_2 és \mathbf{F}_2 mátrixok a szorzás szerint jól particionáltak.

$$\mathbf{E}_2 = \left[\begin{array}{cc} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{array} \right] = \frac{\begin{array}{ccc|c} a_{11} & a_{12} & a_{13} & b_{11} \\ a_{21} & a_{22} & a_{23} & b_{21} \\ c_{11} & c_{12} & c_{13} & d_{11} \\ c_{21} & c_{22} & c_{23} & d_{21} \end{array}}{\begin{array}{ccc|c} j_{11} & j_{12} & j_{13} & k_{11} \\ j_{21} & j_{22} & j_{23} & k_{21} \\ l_{11} & l_{12} & l_{13} & m_{11} \\ l_{21} & l_{22} & l_{23} & m_{21} \end{array}}, \quad \mathbf{F}_2 = \left[\begin{array}{cc} \mathbf{J} & \mathbf{K} \\ \mathbf{L} & \mathbf{M} \end{array} \right] = \frac{\begin{array}{ccc|c} j_{11} & j_{12} & j_{13} & k_{11} \\ j_{21} & j_{22} & j_{23} & k_{21} \\ l_{11} & l_{12} & l_{13} & m_{11} \\ l_{21} & l_{22} & l_{23} & m_{21} \end{array}}{\begin{array}{ccc|c} j_{11} & j_{12} & j_{13} & k_{11} \\ j_{21} & j_{22} & j_{23} & k_{21} \\ l_{11} & l_{12} & l_{13} & m_{11} \\ l_{21} & l_{22} & l_{23} & m_{21} \end{array}}$$

Azonban az összeadás szerint \mathbf{E}_2 és \mathbf{F}_2 nem jól particionáltak, mert a megfelelő minormátrixok összege nem képezhető.

Definíció. Legyen \mathbf{A} $n \times n$ -es mátrix. Legyen \mathbf{A}_{11} az \mathbf{A} mátrix egy minormátrixa. Ha az \mathbf{A}_{11} minormátrix elemei az \mathbf{A} mátrix főátlójára szimmetrikusan helyezkednek el, az \mathbf{A}_{11} minormátrixot főminormátrixnak nevezzük.

Definíció. Legyen \mathbf{A} $n \times n$ -es mátrix. Legyen \mathbf{A} egy particionálása

$$\mathbf{A} = \left[\begin{array}{cc} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{array} \right],$$

ahol \mathbf{A}_{11} $r \times r$ -es ($r < n$) főminormátrix. Ekkor \mathbf{A}_{11} az \mathbf{A} mátrix bal felső r -ed rendű sarokminormátrixa.

Példa. Legyen

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix}.$$

Az \mathbf{A} mátrix bal felső elsőrendű sarokminormátrixa $[a_{11}]$, a bal felső másodrendű sarokminormátrixa

$$\begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix},$$

és a bal felső harmadrendű sarokminormátrixa önmaga.

Az alkalmazásokban gyakran szükség van komplex vektorokra, mátrixokra. A belső szorzatot leszámítva a fenti definíciók érvényesek komplex vektorokra és mátrixokra is, azzal a különbséggel, hogy a komplex vektorok elemei komplex számok, a komplex mátrixok elemei komplex vektorok. A belső szorzat általánosabb definíciója megköveteli, hogy egy vektor saját magával vett belső szorzata valós, nem negatív, és csak akkor nulla, ha a vektor nulla. Defináljuk a belső szorzatot komplex vektorokra is.

Definíció. A $\mathbf{z} = \mathbf{x} + i\mathbf{y}$ komplex vektor konjugáltja $\bar{\mathbf{z}} = \mathbf{x} - i\mathbf{y}$.

Definíció. Az $\mathbf{A} = [a_{ij}]$ komplex mátrix konjugáltja $\bar{\mathbf{A}} = [\bar{a}_{ij}]$.

Definíció. A $\mathbf{z} = \mathbf{x} + i\mathbf{y}$ komplex vektor Hermite-féle transzponáltja

$$\mathbf{z}^H = \bar{\mathbf{z}}^T = \mathbf{x}^T - i\mathbf{y}^T. \quad (3.5)$$

Definíció. Az $\mathbf{A} = \mathbf{B} + i\mathbf{C}$ komplex mátrix Hermite-féle transzponáltja

$$\mathbf{A}^H = \bar{\mathbf{A}}^T = \mathbf{B}^T - i\mathbf{C}^T. \quad (3.6)$$

Definíció. Az \mathbf{A} komplex mátrix Hermite-mátrix, ha $\mathbf{A} = \mathbf{A}^H$.

Példa. Legyen $\mathbf{A} = \mathbf{B} + i\mathbf{C}$ és $\mathbf{z} = \mathbf{x} + i\mathbf{y}$. Ekkor $\mathbf{A}^H \mathbf{z} = (\mathbf{B}^T - i\mathbf{C}^T)(\mathbf{x} + i\mathbf{y}) = \mathbf{B}^T \mathbf{x} + \mathbf{C}^T \mathbf{y} + i(\mathbf{B}^T \mathbf{y} - \mathbf{C}^T \mathbf{x})$.

A komplex Hermite-mátrixot így a valós szimmetrikus mátrix analógiájára definiáltuk. Hasonlóan a valós szimmetrikus esethez, a komplex esetre is igaz, hogy $(\mathbf{A}^H \mathbf{A})^H = \mathbf{A}^H \mathbf{A}$.

Definíció. A \mathbf{w} komplex vektor és a \mathbf{z} komplex vektor belső szorzata $\mathbf{z}^H \mathbf{w}$.

A $\mathbf{z}^H \mathbf{z} = (\mathbf{x}^T - i\mathbf{y}^T)(\mathbf{x} + i\mathbf{y}) = \mathbf{x}^T \mathbf{x} + \mathbf{y}^T \mathbf{y}$ szorzat nem lehet se komplex, se negatív, és csak akkor nulla, ha \mathbf{z} vektor nulla, így a definíció eleget tesz az általánosabb feltételeknek.

3.2. Lineáris függetlenség

A lineáris függetlenség alapvető fogalom. A következő fejezetekben nagyban fogunk annak a következményeire támaszkodni, hogy vektorok lineárisan függetlenek-e, vagy sem.

Definíció. Az \mathbf{a}_i , $i = 1, 2, \dots, n$ vektorok lineárisan összefüggők, ha léteznek olyan x_i skalárok, amelyek nem mindegyike zérus, és a velük képzett lineáris kombinációra fennáll, hogy

$$\sum_{i=1}^n \mathbf{a}_i x_i = \mathbf{0}. \quad (3.7)$$

Ellenkező esetben az \mathbf{a}_i vektorok lineárisan függetlenek.

Példa. Az $\mathbf{a}_1^T = [1 \ 2 \ -1]$, $\mathbf{a}_2^T = [-2 \ -1 \ 1]$, $\mathbf{a}_3^T = [-1 \ 4 \ -1]$ vektorok lineárisan összefüggők, mert $3\mathbf{a}_1 + 2\mathbf{a}_2 - \mathbf{a}_3 = \mathbf{0}$.

Példa. Az $\mathbf{a}_1^T = [1 \ 0 \ 0]$, $\mathbf{a}_2^T = [1 \ 1 \ 0]$, $\mathbf{a}_3^T = [-1 \ 1 \ 1]$ vektorok lineárisan függetlenek, mert

$$\sum_{i=1}^n \mathbf{a}_i x_i = \begin{bmatrix} x_1 + x_2 - x_3 \\ x_2 + x_3 \\ x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}.$$

Ebből következik, hogy $x_3 = 0$. Így $x_2 + x_3 = 0$ és $x_1 + x_2 + x_3 = 0$ akkor és csak akkor, ha x_1 és x_2 is nulla.

Definíció. Az \mathbf{A} $m \times n$ -es mátrix oszlopai lineárisan összefüggők, ha létezik olyan n -ed rendű $\mathbf{x} \neq \mathbf{0}$ vektor, hogy $\mathbf{A}\mathbf{x} = \mathbf{0}$. Ha nem létezik ilyen vektor, akkor a mátrix oszlopai lineárisan függetlenek.

Ha az \mathbf{A} mátrix oszlopai lineárisan összefüggők, abból nem csak az következik, hogy az oszlopok megfelelő lineáris kombinációja nulla, hanem hogy létezik olyan $\mathbf{x} \neq \mathbf{0}$ vektor, amely az \mathbf{A} mátrix minden sorára ortogonális. Az \mathbf{x} vektor nem egyértelmű, mert a skalárral szorzása nem változtat az ortogonalitáson, így \mathbf{x} tetszőlegesen méretezhető. Az \mathbf{A} mátrix oszlopainak lineáris függetlenségét nem befolyásolja, ha az \mathbf{A} mátrix sorait felcseréljük.

Definíció. A lineárisan független n -ed rendű vektorok halmazát, amelyből lineáris kombinációként bármely más n -ed rendű vektor kifejezhető, bázisnak nevezzük.

3.1. Tétel. Az $n + 1$ darab n -ed rendű vektor lineárisan összefüggő. [3, 29. oldal, Theorem 2.2]

Bizonyítás. Indukcióval bizonyítunk. Megmutatjuk, hogy ha bármely $(r-1) \times r$ -es mátrix oszlopai lineárisan összefüggők, akkor bármely $r \times (r+1)$ -es mátrix oszlopai is lineárisan összefüggők. Legyen \mathbf{A}_1 $r \times (r+1)$ -es mátrix és

$$\mathbf{A}_1 = \begin{bmatrix} \mathbf{A} & \mathbf{b} \\ \mathbf{c}^T & \delta \end{bmatrix},$$

ahol \mathbf{A} $(r-1) \times r$ -es mátrix, \mathbf{b} oszlopvektor, \mathbf{c}^T sorvektor és δ skalár. Az $\mathbf{A} - \mathbf{b}\delta^{-1}\mathbf{c}^T$ egy $(r-1) \times r$ -es mátrix, és az indukciós feltétel szerint bármely $(r-1) \times r$ -es mátrix oszlopai lineárisan összefüggők. Ezért létezik egy olyan $\mathbf{x} \neq \mathbf{0}$ vektor, hogy

$$(\mathbf{A} - \mathbf{b}\delta^{-1}\mathbf{c}^T)\mathbf{x} = \mathbf{0}.$$

Ebből következik, hogy

$$\begin{bmatrix} \mathbf{A} & \mathbf{b} \\ \mathbf{c}^T & \delta \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \delta^{-1}\mathbf{c}^T \end{bmatrix} = \begin{bmatrix} \mathbf{0} \\ 0 \end{bmatrix}.$$

Mivel $\mathbf{x} \neq \mathbf{0}$, léteznie kell egy olyan $\mathbf{y} \neq \mathbf{0}$ vektornak, hogy $\mathbf{A}_1\mathbf{y} = \mathbf{0}$. Tehát \mathbf{A}_1 sorai is lineárisan összefüggők. \square

Következmény. Az n -ed rendű vektorok bázisa mindig n darab vektor.

3.2. Tétel. Legyen r darab lineárisan független vektorunk, amik egy új vektor hozzáadásával lineárisan összefüggővé válnak. Ekkor az új vektor kifejezhető az eredeti r vektorok egyértelmű lineáris kombinációjaként. [3, 30. oldal, Lemma 2.2]

Bizonyítás. Legyen \mathbf{A} $n \times r$ -es mátrix és \mathbf{b} n -ed rendű vektor. Legyenek az \mathbf{A} mátrix oszlopai lineárisan függetlenek. Legyen a \mathbf{b} vektor az \mathbf{A} mátrix oszlopaival lineárisan összefüggő. Mivel ez összesen $r+1$ darab lineárisan összefüggő vektor, létezik olyan r -ed rendű \mathbf{y} vektor és η skalár, hogy $\mathbf{A}\mathbf{y} + \mathbf{b}\eta = \mathbf{0}$. A skalár η nem lehet nulla, mert akkor \mathbf{A} mátrix oszlopai lineárisan összefüggők lennének, ami ellentmond a hipotézisnek. Így az $\mathbf{y}\eta^{-1}$ vektor létezik, és ha $\mathbf{x} = \mathbf{y}\eta^{-1}$, akkor $\mathbf{b} = \mathbf{A}\mathbf{x}$. Azt kell belátni, hogy az \mathbf{x} vektor egyértelmű. Tegyük fel, hogy létezik olyan \mathbf{z} vektor, hogy $\mathbf{b} = \mathbf{A}\mathbf{z}$. Ha ezt kivonjuk az előző egyenlőségből, azt kapjuk, hogy $\mathbf{0} = \mathbf{A}(\mathbf{x} - \mathbf{z})$. Mivel \mathbf{A} mátrix oszlopai lineárisan függetlenek, ez csak akkor lehetséges, ha $\mathbf{x} = \mathbf{z}$. \square

Definíció. Az \mathbf{A} mátrix lineárisan független oszlopainak maximum darabszáma az \mathbf{A} mátrix rangja. Jelölése $r(\mathbf{A})$. Az \mathbf{A} mátrix teljesrangú, ha rangja megegyezik az oszlopainak számával.

3.3. Nemszinguláris mátrixok

Definiáljuk a nemszinguláris mátrixot. Ehhez megvizsgáljuk a mátrix oszlopainak illetve sorainak lineáris függőségét, definiáljuk a mátrix inverzét, és megvizsgáljuk a lineáris függőség és a mátrix inverze közötti kapcsolatot.

3.3. Tétel. *Nem létezik olyan \mathbf{X} mátrix, hogy $\mathbf{AX} = \mathbf{I}$, ha \mathbf{A} mátrix sorai lineárisan összefüggők. [3, 30. oldal, Lemma 2.3]*

Bizonyítás. Tegyük fel az ellenkezőjét. Mivel \mathbf{A} mátrix lineárisan összefüggő, létezik olyan $\mathbf{y} \neq \mathbf{0}$ vektor, hogy $\mathbf{y}^T \mathbf{A} = \mathbf{0}^T$. Így $\mathbf{y}^T \mathbf{AX} = \mathbf{0}^T$, de mivel $\mathbf{AX} = \mathbf{I}$, ez csak akkor teljesül, ha $\mathbf{y} = \mathbf{0}$. Így ellentmondáshoz jutottunk. \square

Definíció. Legyen \mathbf{A} négyzetes mátrix. Ha létezik olyan \mathbf{X} mátrix, hogy

$$\mathbf{AX} = \mathbf{I}, \quad (3.8)$$

akkor \mathbf{X} az \mathbf{A} mátrix jobboldali inverze. Ha létezik olyan \mathbf{Y} mátrix, hogy

$$\mathbf{YA} = \mathbf{I}, \quad (3.9)$$

akkor \mathbf{Y} az \mathbf{A} mátrix baloldali inverze.

Definíció. Legyen \mathbf{A} négyzetes mátrix. Ha létezik olyan \mathbf{X} mátrix, hogy

$$\mathbf{AX} = \mathbf{XA} = \mathbf{I}, \quad (3.10)$$

akkor \mathbf{X} az \mathbf{A} mátrix inverze. Az \mathbf{A} mátrix inverzének jelölése \mathbf{A}^{-1} . Egy mátrix invertálható, ha létezik inverze.

3.4. Tétel. *Legyen \mathbf{A} n -ed rendű négyzetes mátrix. Ha az \mathbf{A} mátrix oszlopai lineárisan függetlenek, létezik olyan \mathbf{X} mátrix, hogy $\mathbf{AX} = \mathbf{I}$. [3, 30. oldal, Theorem 2.3]*

Bizonyítás. Legyen \mathbf{A} $n \times n$ -es négyzetes mátrix, és legyenek az oszlopai lineárisan függetlenek. Legyenek \mathbf{A}_i ($i = 1, 2, \dots, n-1$) $(n-1) \times n$ -es mátrixok. Az \mathbf{A}_i mátrix mindig legyen az \mathbf{A} mátrix i -edik sorának elhagyásával kapott mátrix. Azaz \mathbf{A}_i j -edik sora az \mathbf{A} mátrix j -edik sora, ha $1 \leq j \leq i-1$, de a $(j+1)$ -edik sora, ha $i \leq j \leq n-1$. A 3.1-es tételből következik, hogy mindig létezik olyan $\mathbf{x}_i \neq \mathbf{0}$, hogy $\mathbf{A}_i \mathbf{x}_i = \mathbf{0}$. Mivel létezik olyan \mathbf{x}_i , hogy $\mathbf{A}_i \mathbf{x}_i = \mathbf{0}$, de $\mathbf{Ax}_i \neq \mathbf{0}$, így az \mathbf{Ax}_i vektornak egyedül az i -edik eleme nem nulla. Minden i -re választhatjuk \mathbf{x}_i -t úgy, hogy $\mathbf{Ax}_i = \mathbf{e}_i$ legyen, ahol \mathbf{e}_i az n -ed rendű egységmátrix i -edik oszlopa. Legyen $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]$, így $\mathbf{AX} = \mathbf{I}$. \square

Megmutattuk, hogy ha az \mathbf{A} négyzetes mátrixnak az oszlopai lineárisan függetlenek, akkor létezik jobboldali inverze. A 3.3-as tételből következik, hogy ekkor az \mathbf{A} mátrix sorai is lineárisan függetlenek. Tehát a négyzetes mátrixok sorainak és oszlopainak lineáris függetlensége egyenértékű.

3.5. Tétel. *Ha egy négyzetes mátrix oszlopai lineárisan függetlenek, akkor jobboldali inverze egyértelmű. [3, 31. oldal, Theorem 2.4]*

Bizonyítás. Legyen \mathbf{A} négyzetes mátrix, és legyenek az oszlopai lineárisan függetlenek. Tegyük fel, hogy az \mathbf{A} mátrix jobboldali inverze nem egyértelmű, azaz léteznek olyan $\mathbf{X} \neq \mathbf{Y}$ mátrixok, hogy $\mathbf{AX} = \mathbf{I}$ és $\mathbf{AY} = \mathbf{I}$. Vonjuk ki a két egyenlőséget egymásból, így kapjuk, hogy $\mathbf{A}(\mathbf{X} - \mathbf{Y}) = \mathbf{0}$. Ez azonban az \mathbf{A} mátrix oszlopainak lineáris függetlensége miatt csak akkor lehetséges, ha $\mathbf{X} = \mathbf{Y}$. \square

3.6. Tétel. *Ha egy négyzetes mátrixnak létezik egyértelmű jobboldali inverze, akkor az megegyezik a baloldali inverzével. [3, 31. oldal, Theorem 2.5]*

Bizonyítás. Legyen \mathbf{A} négyzetes mátrix, és legyenek az oszlopai lineárisan függetlenek. Legyen $\mathbf{AX} = \mathbf{I}$. Az egyenlőséget jobbról \mathbf{A} mátrixszal szorozva kapjuk, hogy $\mathbf{AXA} = \mathbf{A}$, átrendezve $\mathbf{A}(\mathbf{XA} - \mathbf{I}) = \mathbf{0}$. Ebből az \mathbf{A} mátrix oszlopainak lineáris függetlensége miatt következik, hogy $\mathbf{XA} = \mathbf{I}$. \square

Definíció. Legyen \mathbf{A} négyzetes mátrix. Ha a következő ekvivalens állítások teljesülnek, akkor az \mathbf{A} mátrix nonszinguláris, egyébként szinguláris.

1. Az \mathbf{A} mátrix oszlopai lineárisan függetlenek.
2. Az \mathbf{A} mátrix sorai lineárisan függetlenek.
3. Az \mathbf{A} mátrix invertálható.

3.7. Tétel. *Legyen \mathbf{A} négyzetes mátrix, és legyen*

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{0} & \mathbf{A}_{22} \end{bmatrix},$$

ahol \mathbf{A}_{11} és \mathbf{A}_{22} minormátrixok is négyzetes mátrixok. Ekkor az \mathbf{A} mátrix akkor és csak akkor nonszinguláris, ha \mathbf{A}_{11} és \mathbf{A}_{22} minormátrixok nonszingulárisak. [3, 33. oldal, Lemma 2.4]

Bizonyítás. Legyenek \mathbf{A}_{11} és \mathbf{A}_{22} minormátrixok nonszingulárisak. Ekkor

$$\mathbf{A}^{-1} = \begin{bmatrix} \mathbf{A}_{11}^{-1} & \mathbf{A}_{11}^{-1}\mathbf{A}_{12}\mathbf{A}_{22}^{-1} \\ \mathbf{0} & \mathbf{A}_{22}^{-1} \end{bmatrix},$$

tehát az \mathbf{A} mátrix nemszinguláris.

Legyen \mathbf{A}_{11} minormátrix szinguláris. Ekkor mindig létezik olyan $\mathbf{x} \neq \mathbf{0}$ vektor, hogy $\mathbf{A}_{11}\mathbf{x} = \mathbf{0}$. Ekkor az \mathbf{A} mátrix szingularitása csak \mathbf{A}_{11} minormátrixtól függ, mert

$$\begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{0} & \mathbf{A}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{0} \end{bmatrix} = \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}.$$

Ezért, ha \mathbf{A}_{11} minormátrix szinguláris, akkor az \mathbf{A} mátrix szinguláris. Hasonlóan, ha \mathbf{A}_{22} minormátrix szinguláris, mindig létezik olyan $\mathbf{y} \neq \mathbf{0}$ vektor, hogy $\mathbf{y}^T \mathbf{A}_{22} = \mathbf{0}^T$. \square

Most definiálunk és megvizsgálunk néhány elméleti vagy gyakorlati szempontból fontos nemszinguláris mátrixot.

Definíció. Az ortogonális mátrix olyan valós mátrix, melynek az inverze a transzponáltja.

Definíció. Az $\mathbf{U} = [u_{ij}]$ négyzetes mátrix felső háromszögmátrix, ha $u_{ij} = 0$ minden $i > j$ -re.

Definíció. Az $\mathbf{L} = [l_{ij}]$ négyzetes mátrix alsó háromszögmátrix, ha $l_{ij} = 0$ minden $i < j$ -re.

Definíció. Az alsó háromszögmátrix alsó egység háromszögmátrix, ha minden diagonális eleme egy.

Definíció. A felső háromszögmátrix felső egység háromszögmátrix, ha minden diagonális eleme egy.

3.8. Tétel. Az $\mathbf{U} = [u_{ij}]$ felső háromszögmátrix akkor és csak akkor nemszinguláris, ha $u_{ii} \neq 0$ minden i -re. [3, 34. oldal, Theorem 2.6]

Bizonyítás. Legyen \mathbf{U} n -ed rendű felső háromszögmátrix. Legyen \mathbf{U}_i a bal felső i -ed rendű sarokminormátrixa \mathbf{U} -nak., azaz $\mathbf{U} = \mathbf{U}_n$. Legyen

$$\mathbf{U} = \begin{bmatrix} \mathbf{U}_{i-1} & \mathbf{v}_i \\ \mathbf{0} & u_{ii} \end{bmatrix}, \quad 2 \leq i \leq n$$

és $\mathbf{v}_i^T = [u_{1i}, u_{2i}, \dots, u_{(i-1)i}]$. A 3.7-es tételből következik, hogy \mathbf{U} akkor és csak akkor nemszinguláris, ha minden \mathbf{U}_i sarokminormátrix nemszinguláris. Ha $u_{ii} \neq 0$ bármely $1 \leq i \leq n$ esetén, akkor $\mathbf{U}_1, \mathbf{U}_2, \dots, \mathbf{U}_n$ sarokminormátrixok nemszingulárisak, azaz \mathbf{U} nemszinguláris. \square

3.9. Tétel. Négyzetes mátrixok szorzata akkor és csak akkor nemszinguláris, ha a szorzat minden tényezője nemszinguláris.

Bizonyítás. Legyenek \mathbf{A} , \mathbf{B} és \mathbf{C} négyzetes mátrixok. Legyen $\mathbf{C} = \mathbf{AB}$. Tegyük fel, hogy \mathbf{B} és \mathbf{C} mátrixok nemszingulárisak, de \mathbf{A} mátrix szinguláris. Mivel \mathbf{C} nemszinguláris, létezik inverze, és $\mathbf{CC}^{-1} = \mathbf{I}$. Ebből következik, hogy $(\mathbf{AB})\mathbf{C}^{-1} = \mathbf{I}$. De ez csak akkor lehetséges, ha $\mathbf{C}^{-1} = \mathbf{B}^{-1}\mathbf{A}^{-1}$, amivel ellentmondáshoz jutottunk. \square

Definíció. Az \mathbf{A} valós mátrix pozitív definit, ha $\mathbf{x}^T \mathbf{A} \mathbf{x} > 0$, minden $\mathbf{x} \neq \mathbf{0}$ vektorra. Jelölése $\mathbf{A} > 0$.

Definíció. Az \mathbf{A} valós mátrix pozitív szemidefinit, ha $\mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0$, minden $\mathbf{x} \neq \mathbf{0}$ vektorra. Jelölése $\mathbf{A} \geq 0$.

Definíció. Az \mathbf{A} valós mátrix negatív definit, ha $\mathbf{x}^T \mathbf{A} \mathbf{x} < 0$, minden $\mathbf{x} \neq \mathbf{0}$ vektorra. Jelölése $\mathbf{A} < 0$.

Definíció. Az \mathbf{A} valós mátrix negatív szemidefinit, ha $\mathbf{x}^T \mathbf{A} \mathbf{x} \leq 0$, minden $\mathbf{x} \neq \mathbf{0}$ vektorra. Jelölése $\mathbf{A} \leq 0$.

Definíció. Az \mathbf{A} valós mátrix definit, ha $\mathbf{x}^T \mathbf{A} \mathbf{x} \neq 0$, minden $\mathbf{x} \neq \mathbf{0}$ vektorra.

Definíció. Az \mathbf{A} valós mátrix indefinit, ha léteznek olyan $\mathbf{x} \neq \mathbf{0}$ és $\mathbf{y} \neq \mathbf{0}$ vektorok, amelyekre $\mathbf{x}^T \mathbf{A} \mathbf{x} < 0 < \mathbf{y}^T \mathbf{A} \mathbf{y}$.

3.10. Tétel. Legyen \mathbf{A} $m \times n$ -es valós mátrix lineárisan független oszlopokkal. Ekkor $\mathbf{A}^T \mathbf{A}$ pozitív definit. [3, 34. oldal, Lemma 2.5]

Bizonyítás. Legyen \mathbf{A} $m \times n$ -es valós mátrix lineárisan független oszlopokkal. Az $\mathbf{A}^T \mathbf{A}$ mátrix szimmetrikus. Legyen $\mathbf{y} = \mathbf{Ax}$, így $\mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{y}^T \mathbf{y} > 0$, akkor és csak akkor, ha $\mathbf{y} \neq \mathbf{0}$. Mivel az \mathbf{A} mátrix oszlopai lineárisan függetlenek, \mathbf{y} akkor és csak akkor nulla, ha \mathbf{x} is nulla. Az $\mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x}$ akkor és csak akkor nulla, ha $\mathbf{x} = \mathbf{0}$. \square

3.4. Vektornormák és mátrixnormák

Gyakran akarunk vektorokat vagy mátrixokat a nagyságuk alapján összehasonlítani. Például ha egy iteratív eljárással közelítünk egy vektorhoz, a hibavektor (azaz a vektor és a közelítő vektor közötti különbség) nagysága jó ha gyorsan csökken. A norma a vektorokhoz és a mátrixokhoz egy skalárt rendel.

Definíció. Az $\|\mathbf{x}\|$ skalár az \mathbf{x} vektor normája, ha kielégíti a következő három feltételt.

1. $\|\mathbf{x}\| = 0$, ha $\mathbf{x} = \mathbf{0}$, egyébként $\|\mathbf{x}\| > 0$.
2. $\|\mathbf{x}\theta\| = \|\mathbf{x}\| |\theta|$, ahol θ skalár.

$$3. \quad \|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|.$$

Definíció. Az \mathbf{x} vektor l_1 , l_2 , l_∞ normáinak definíciója:

1. $\|\mathbf{x}\|_1 = \sum_i |x_i|$ az l_1 norma,
2. $\|\mathbf{x}\|_2 = (\sum_i |x_i|^2)^{\frac{1}{2}}$ az l_2 (euklideszi) norma,
3. $\|\mathbf{x}\|_\infty = \max_i |x_i|$ az l_∞ (maximum) norma.

3.11. Tétel. (Cauchy-egyenlőtlenség) Legyen \mathbf{x} és \mathbf{y} n -ed rendű nem nulla vektor. Az

$$|\mathbf{x}^T \mathbf{y}| \leq \|\mathbf{x}\|_2 \|\mathbf{y}\|_2 \quad (3.11)$$

egyenlőtlenség akkor és csak akkor igaz, ha \mathbf{y} vektor skalárszorosa \mathbf{x} vektornak. [3, 42. oldal, Cauchy's Inequality]

A Cauchy-egyenlőtlenség komplex vektorokra is igaz, ha a transzponáltat Hermite-féle transzponáltra cseréjük. A tételt itt nem bizonyítjuk.

Definíció. Az $\|\mathbf{A}\|$ skalár az \mathbf{A} mátrix normája, ha kielégíti a következő négy feltételt.

1. $\|\mathbf{A}\| = 0$, ha $\mathbf{A} = \mathbf{0}$, egyébként $\|\mathbf{A}\| > 0$.
2. $\|\mathbf{A}\theta\| = \|\mathbf{A}\| |\theta|$, ahol θ skalár.
3. $\|\mathbf{A} + \mathbf{B}\| \leq \|\mathbf{A}\| + \|\mathbf{B}\|$.
4. $\|\mathbf{AB}\| \leq \|\mathbf{A}\| \|\mathbf{B}\|$.

Definíció. Az \mathbf{A} mátrix $\|\mathbf{A}\|_p$ indukált normájának definíciója

$$\|\mathbf{A}\|_p = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{Ax}\|_p}{\|\mathbf{x}\|_p}, \quad (3.12)$$

ahol $p = 1, 2$ vagy ∞ .

Az indukált mátrixnormákra mindig igaz, hogy $\|\mathbf{Ax}\| \leq \|\mathbf{A}\| \|\mathbf{x}\|$. Ez a gyakorlatban sokszor hasznos. Az l_1 indukált mátrixnormát explicit az

$$\|\mathbf{A}\|_1 = \max_j \sum_i |a_{ij}| \quad (3.13)$$

képlettel számolhatjuk. Az l_∞ indukált mátrixnormát az

$$\|\mathbf{A}\|_\infty = \max_i \sum_j |a_{ij}| \quad (3.14)$$

képlettel számolhatjuk. Az l_2 által indukált mátrixnormát spektrális normának hívjuk. A spektrális norma elméleti fontossága mellett nagy gyakorlati hátránya, hogy nincs egyszerű explicit képlet a kiszámolására.

Definíció. Az \mathbf{A} mátrix Frobenius normájának definíciója

$$\|\mathbf{A}\|_F = \left(\sum_{i,j} |a_{ij}|^2 \right)^{\frac{1}{2}}. \quad (3.15)$$

A Frobenius norma az euklideszi vektornorma mátrix megfelelője, de nem az euklideszi vektornorma által indukált mátrixnorma. Nem vektornorma által indukált mátrixnorma, ezért a gyakorlatban sokszor indokolatlanul pontatlan eredményekhez vezet.

Definíció. A nonszinguláris \mathbf{A} mátrix kondíciós számának jelölése $k(\mathbf{A})$, és definíciója

$$k(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^{-1}\|, \quad (3.16)$$

ahol tetszőleges norma választható.

A kondíciós szám mérete jellemzi a mátrix szingularitásának mértékét. Ha a kondíciós szám kicsi, szokás a mátrixot jól kondicionáltnak hívni, ha a kondíciós szám nagy, szokás a mátrixot rosszul kondicionáltnak, vagy közel szingulárisnak hívni [3]. Ezt bővebben lásd a 8.2. fejezetben.

4. Lineáris egyenletrendszerek megoldása, osztályozása

A megoldása az

$$\mathbf{A}\mathbf{x} = \mathbf{b} \quad (4.1)$$

alakú egyenletrendszernek, ahol \mathbf{A} n -ed rendű nonszinguláris mátrix, \mathbf{x} ismeretlen n -ed rendű vektor, \mathbf{b} tetszőleges n -ed rendű vektor, az egyik leggyakoribb feladat. Az \mathbf{A} mátrix az egyenletrendszer együtthatómátrixa. A (4.1) egyenletrendszernek akkor és csak akkor létezik megoldása, ha a \mathbf{b} vektor lineárisan kifejezhető az \mathbf{A} mátrix oszlopvektoraiból. Vagy másként fogalmazva, a \mathbf{b} vektor és az \mathbf{A} mátrix oszlopvektorai nem lineárisan függetlenek. Ekkor az egyenletrendszer megoldása

$$\mathbf{A}^{-1}\mathbf{b}. \quad (4.2)$$

Az \mathbf{A} mátrix inverzének számolása a gyakorlatban nem célszerű, mert túl műveletigényes. A megoldási módszerek két fő osztálya az iteratív eljárások és a direkt

eljárások. Egy megoldási módszer hatékonysága két fő szempont alapján bírálható el:

1. A megoldási módszer mennyire gyors, azaz mennyire műveletigényes?
2. Mennyire pontos a kiszámított megoldás? [14]

A gyakorlatban előforduló együtthatómátrixok általában vagy kitöltöttek és alacsony rendszámúak (pl. a rendszám kisebb mint 30 [14]), vagy ritkák és nagy rendszámúak. A direkt eljárások általában előnyösebbek a kis rendszámú kitöltött mátrixoknál, míg az iteratív eljárások általában előnyösebbek a nagy rendszámú ritka mátrixoknál. [14]

Az egyenletrendszer megoldásánál számolnunk kell különböző hibaforrásokkal. Hibaforrás lehet az együtthatómátrix vagy a jobboldali \mathbf{b} vektor elemeinek mérési hibája, a számítások közben keletkező kerekítési hibák, és a megoldási módszer képlethibája. [14]

4.1. Direkt eljárások

A direkt eljárások elvileg mentesek a képlethibától, és a megoldást véges sok lépésben állítják elő. Elméletben a megoldás pontos, a gyakorlatban azonban ez mégsem teljesül, a lépések közben keletkező kerekítési hibák miatt. A direkt eljárások alapja az (4.1) egyenletrendszer sorozatos transzformációja, amíg a megoldás könnyen számolhatóvá válik.

4.1.1. Az LU-felbontás

Definíció. Az \mathbf{A} mátrix LU-felbontásán az

$$\mathbf{A} = \mathbf{L}\mathbf{U} \quad (4.3)$$

felbontást értjük, ahol \mathbf{L} alsó háromszögmátrix és \mathbf{U} felső háromszögmátrix.

Az LU-felbontással az (4.1) egyenletrendszer megoldása az

$$\mathbf{L}\mathbf{y} = \mathbf{b} \quad (4.4)$$

és

$$\mathbf{U}\mathbf{x} = \mathbf{y} \quad (4.5)$$

háromszög alakú egyenletrendszerek megoldására egyszerűsíthető, ahol \mathbf{y} vektor a felbontás egy köztes eredménye. Először a (4.5) egyenletrendszert oldjuk meg. Legyen

$\mathbf{U} = [u_{ij}]$, $\mathbf{x} = [x_i]$ és $\mathbf{y} = [y_i]$ ($i, j = 1, 2, \dots, n$) ekkor

$$x_i = (y_i - \sum_{j=i+1}^n u_{ij}x_j)/u_{ii}, \quad (4.6)$$

ahol $u_{ii} \neq 0$. Ahhoz, hogy a (4.5) egyenletrendszernek tetszőleges \mathbf{y} vektorra legyen megoldása, feltételezzük, hogy \mathbf{U} mátrix nonszinguláris. Így a 3.8-es tételből következik, hogy $u_{ii} \neq 0$. A (4.4) egyenletrendszer \mathbf{y} vektor ismeretében hasonlóan oldható meg. [3]

4.1. Tétel. *Az \mathbf{A} nonszinguláris mátrixnak akkor és csak akkor létezik \mathbf{LU} -felbontása, ha az \mathbf{A} mátrix minden bal felső sarokminormátrixa nonszinguláris. [3, 56. oldal, Theorem 4.1]*

Bizonyítás. Legyen \mathbf{A} mátrix nonszinguláris. Megmutatjuk, hogy ha az $\mathbf{A} = \mathbf{LU}$ felbontás létezik, akkor az \mathbf{A} mátrix bal felső sarokminormátrixai szükségképpen nonszingulárisak. Particionáljuk \mathbf{A} mátrixot úgy, hogy

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix},$$

ahol \mathbf{A}_{11} egy tetszőleges rendű bal felső sarokminormátrixa az \mathbf{A} mátrixnak. Particionáljuk az \mathbf{L} és \mathbf{U} mátrixokat hasonlóan. A hipotézis szerint $\mathbf{A} = \mathbf{LU}$, azaz

$$\begin{bmatrix} \mathbf{L}_{11} & \mathbf{0} \\ \mathbf{L}_{21} & \mathbf{L}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{U}_{11} & \mathbf{U}_{12} \\ \mathbf{0} & \mathbf{U}_{22} \end{bmatrix} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix}$$

és

$$\mathbf{L}_{11}\mathbf{U}_{11} = \mathbf{A}_{11}.$$

Az \mathbf{A} mátrix nonszinguláris, így a 3.9-es tétel miatt \mathbf{L} és \mathbf{U} mátrixok is nonszingulárisak. A 3.7-es tétel miatt \mathbf{L}_{11} és \mathbf{U}_{11} szintén nonszingulárisak, így \mathbf{A}_{11} nonszinguláris. Mivel \mathbf{A}_{11} egy tetszőleges rendű bal felső sarokminormátrixa az \mathbf{A} mátrixnak, így az \mathbf{A} mátrix szükségképpen nonszinguláris.

Legyen

$$\mathbf{A}_r = \mathbf{L}_r\mathbf{U}_r, \quad (4.7)$$

ahol $\mathbf{A}_r, \mathbf{L}_r$ és \mathbf{U}_r az \mathbf{A}, \mathbf{L} és \mathbf{U} mátrixok r -ed rendű bal felső sarokminormátrixai. Megmutatjuk, hogy ha \mathbf{A}_r nonszinguláris, akkor \mathbf{A}_{r+1} is szükségképpen nonszinguláris. Indukcióval bizonyítunk. Particionáljuk \mathbf{A}_{r+1} -t úgy, hogy

$$\mathbf{A}_{r+1} = \begin{bmatrix} \mathbf{A}_r & \mathbf{b}_r \\ \mathbf{c}_r^T & \delta_r \end{bmatrix}, \quad (4.8)$$

ahol \mathbf{b}_r és \mathbf{c}_r r -ed rendű vektorok és δ_r skalár. Particionáljuk \mathbf{L}_{r+1} -et és \mathbf{U}_{r+1} -et úgy, hogy

$$\mathbf{L}_{r+1} = \begin{bmatrix} \mathbf{L}_r & \mathbf{0} \\ \mathbf{p}_r & \lambda_r \end{bmatrix} \quad (4.9)$$

és

$$\mathbf{U}_{r+1} = \begin{bmatrix} \mathbf{U}_r & \mathbf{v}_r \\ \mathbf{0}^T & \mu_r \end{bmatrix}. \quad (4.10)$$

Ekkor található olyan $\mathbf{p}_r, \lambda_r, \mathbf{v}_r$ és μ_r , hogy

$$\mathbf{L}_{r+1}\mathbf{U}_{r+1} = \mathbf{A}_{r+1}. \quad (4.11)$$

A particionált (4.8), (4.9) és (4.10) mátrixokkal felírva a (4.11) egyenlőséget

$$\begin{bmatrix} \mathbf{L}_r & \mathbf{0} \\ \mathbf{p}_r & \lambda_r \end{bmatrix} \begin{bmatrix} \mathbf{U}_r & \mathbf{v}_r \\ \mathbf{0}^T & \mu_r \end{bmatrix} = \begin{bmatrix} \mathbf{A}_r & \mathbf{b}_r \\ \mathbf{c}_r^T & \delta_r \end{bmatrix}.$$

A szorzás elvégzése után a (4.7) egyenlőséget és az

$$\mathbf{L}_r \mathbf{v}_r = \mathbf{b}_r, \quad (4.12)$$

$$\mathbf{p}_r^T \mathbf{U}_r = \mathbf{c}_r^T \quad (4.13)$$

és

$$\mathbf{p}_r^T \mathbf{v}_r + \lambda_r \mu_r = \delta_r, \quad (4.14)$$

egyenlőségeket kapjuk. Mivel a hipotézis szerint \mathbf{A}_r nonszinguláris, a (4.7) egyenlőségből következik, hogy \mathbf{L}_r és \mathbf{U}_r szintén nonszingulárisak. Ezért léteznek olyan egyértelmű \mathbf{p}_r és \mathbf{v}_r vektorok, amelyek kielégítik a (4.12) és (4.13) egyenlőségeket. Mivel léteznek olyan λ_r és μ_r skalárok, amik kielégítik a (4.14) egyenlőséget, létezik olyan \mathbf{L}_{r+1} és \mathbf{U}_{r+1} , amik kielégítik a (4.11) egyenlőséget. Tehát, ha az \mathbf{A} mátrix bal felső r -ed rendű sarokminormátrixa nonszinguláris, és az \mathbf{A} mátrixnak létezik LU-felbontása, akkor az \mathbf{A} mátrix bal felső $r + 1$ -ed rendű sarokminormátrixa is nonszinguláris. \square

Következmény. Ha \mathbf{L} alsó egység háromszögmátrix, az LU-felbontás egyértelmű. [3, 57. oldal, Corollary]

Bizonyítás. Ha \mathbf{L} alsó egység háromszögmátrix, akkor a (4.9) felbontásban $\lambda_r = 1$ minden r -re. Ebből következik, hogy (4.10) felbontásban μ_r egyértelmű minden r -re. Mivel \mathbf{p}_r és \mathbf{v}_r egyértelmű minden r -re, következik az eredmény. \square

4.1.2. A Gauss-féle elimináció

Az egyik legrégebbi és legjobb módszer a (4.1) egyenletrendszer megoldására a Gauss-féle elimináció. Az első, második, stb. egyenlőség megfelelő skalárszorosát kivonjuk a többi egyenlőségből, úgy, hogy ezzel az ismeretlen változókat elimináljuk. Ezt addig folytatjuk, amíg az utolsó egyenlőségben már csak egy ismeretlen változó marad. Ekkor a kapott egyenlőségek felső háromszög alakúak. A módszert egy példán mutatjuk be.

Példa. Oldjuk meg a

$$\begin{array}{rrrrr} 2x_1 & +x_2 & +3x_3 & -x_4 & = & 3 \\ -4x_1 & -3x_2 & -4x_3 & +5x_4 & = & 2 \\ 6x_1 & +4x_2 & +4x_3 & -5x_4 & = & -1 \\ -4x_1 & -3x_2 & +2x_3 & +4x_4 & = & 7 \end{array}$$

egyenletrendszert! Vonjuk ki az 1. egyenlőség (-2) -szeresét a második és negyedik egyenlőségből, és a 3-szorosát a harmadik egyenlőségből. Így ezekből az egyenlőségekből az x_1 ismeretlent elimináltuk, és a

$$\begin{array}{rrrrr} 2x_1 & +x_2 & +3x_3 & -x_4 & = & 3 \\ & -x_2 & +2x_3 & +3x_4 & = & 8 \\ & & x_2 & -5x_3 & -2x_4 & = & -10 \\ & & -3x_3 & +8x_4 & +2x_4 & = & 13 \end{array}$$

új egyenletrendszert kaptuk. Az új egyenletrendszerben az x_2 ismeretlen eliminálásához vonjuk ki a 2. egyenlőség (-1) -szeresét a harmadik egyenlőségből, és az 1-szeresét a negyedik egyenlőségből. Így ezekből az egyenlőségekből az x_2 ismeretlent elimináltuk, és a

$$\begin{array}{rrrrr} 2x_1 & +x_2 & +3x_3 & -x_4 & = & 3 \\ & -x_2 & +2x_3 & +3x_4 & = & 8 \\ & & -3x_3 & x_4 & = & -2 \\ & & +6x_3 & x_4 & = & 5 \end{array}$$

új egyenletrendszert kaptuk. Végül vonjuk ki a 3. egyenlőség (-2) -szeresét a negyedik egyenlőségből. Így a negyedik egyenlőségből x_2 -t elimináltuk, és a

$$\begin{array}{rrrrr} 2x_1 & +x_2 & +3x_3 & -x_4 & = & 3 \\ & -x_2 & +2x_3 & +3x_4 & = & 8 \\ & & -3x_3 & x_4 & = & -2 \\ & & & x_4 & = & 1 \end{array}$$

új egyenletrendszert kaptuk. Visszahelyettesítve az $x_4 = 1$ -et kapjuk, hogy $x_3 = 1$, $x_2 = -3$ és $x_1 = 2$. [3, 60. oldal, Example 4.3]

Legyen $\mathbf{A}^{(r-1)} = [a_{ij}^{(r-1)}]$, $1 \leq r \leq n$ a (4.1) egyenletrendszer Gauss-féle eliminációval való megoldásakor az $(r-1)$ -edik eliminációs lépés után kapott együtthatómátrix. Legyen $\mathbf{M}_r = [m_{ij}]$ és

$$\mathbf{A}^{(r)} = \mathbf{M}_r \mathbf{A}^{(r-1)}. \quad (4.15)$$

A (4.15) egyenlőségben az \mathbf{M}_r mátrixot mindig úgy választjuk, hogy a szorzás eredményeként az $\mathbf{A}^{(r)}$ együtthatómátrix r -edik oszlopában az r -edik sor alatt az együtthatókra nullákat kapjunk. Ehhez a (4.15) egyenlőségben

$$\mathbf{M}_r = \mathbf{I} - \mathbf{m}_r \mathbf{e}_r^T, \quad (4.16)$$

ahol

$$\mathbf{m}_r = [0, 0, \dots, 0, m_{r+1,r}, m_{r+2,r}, \dots, m_{nr}]^T$$

és

$$m_{ir} = \frac{a_{ir}^{(r-1)}}{a_{rr}^{(r-1)}}, \quad a_{rr}^{(r-1)} \neq 0, \quad (4.17)$$

az \mathbf{e}_r vektor pedig az r -edik oszlopa az \mathbf{I} egységmátrixnak. Ekkor az $\mathbf{A}^{(r)}$ mátrixot úgy számolhatjuk, hogy kivonjuk az $\mathbf{A}^{(r-1)}$ mátrix r -edik sorának m_{ir} -szeresét az $\mathbf{A}^{(r-1)}$ mátrix i -edik sorából minden i -re, ahol $i = r+1, r+2, \dots, n$. [3]

Definíció. A (4.15) egyenlőségben az m_{ir} elemeket multiplikátoroknak hívjuk, az $a_{rr}^{(r-1)}$ elemeket pedig r -edik pivot, vagy r -edik főelemnek hívjuk.

A módszer az alkalmazásakor az

$$\mathbf{A}^{(r)} \mathbf{x} = \mathbf{b}^{(r)} \quad (4.18)$$

egyenlőségrendszereket generálja, ahol $\mathbf{b}^{(r)} = \mathbf{M}_r \mathbf{b}^{(r-1)}$.

$$\mathbf{A}^{(n-1)} = \mathbf{M}_{n-1} \mathbf{M}_{n-2} \dots \mathbf{M}_1 \mathbf{A}, \quad (4.19)$$

ahol $\mathbf{A}^{(n-1)}$ felső háromszögmátrix, és \mathbf{M}_r a (4.16) egyenlőség miatt mindig alsó egység háromszögmátrix. Alsó egység háromszögmátrix inverze alsó egység háromszögmátrix, és alsó egység háromszögmátrixok szorzata is alsó egység háromszögmátrix, így ha

$$\mathbf{U} = \mathbf{A}^{(n-1)} \quad (4.20)$$

és

$$\mathbf{L} = (\mathbf{M}_{n-1} \mathbf{M}_{n-2} \dots \mathbf{M}_1)^{-1}, \quad (4.21)$$

a módszer eredménye az \mathbf{A} együtthatómátrix LU-felbontását adja. Mivel \mathbf{L} alsó egység háromszögmátrix, az LU-felbontás egyértelmű. Ha \mathbf{x} kielégíti a (4.18) egyenlőséget, akkor kielégíti a (4.1) egyenlőséget is. [3]

A (4.17) egyenlőségből látszik, hogy a pivot elemek nem lehetnek nullák. Ha bármelyik pivot elem nulla, akkor az \mathbf{A} mátrix bal felső sarokminormátrixa szinguláris. Ebből következik, hogy az \mathbf{A} mátrixnak nem létezik LU-felbontása (4.1-es tétel). [3]

A nulla értékű r -edik pivot elem kiküszöbölésének egy módja, ha az $\mathbf{A}^{(r-1)}$ mátrix sorait felcseréljük, úgy hogy az r -edik pivot elem helyén ne nullát kapjunk. Ezt úgy érhetjük el, ha az $\mathbf{A}^{(r-1)}$ mátrix r -edik oszlopában az $a_{rr}^{(r-1)}$ elemtől lefelé megkeressük az első nem nulla elemet, legyen ez pl. $a_{ir}^{(r-1)}$, majd felcseréljük az $\mathbf{A}^{(r-1)}$ mátrix r -edik sorát az i -edik sorával. Ezt megtehetjük, mert $i > r$. [3]

4.2. Iteratív eljárások

Az iteratív megoldási módszereknél a (4.1) egyenletrendszer \mathbf{A} mátrixa lényegében változatlan marad és közelítő \mathbf{x}_i ($i = 1, 2, \dots$) megoldásokat generálunk. Azt reméljük, hogy az $\{\mathbf{x}_i\}$ sorozat a megoldáshoz konvergál. A közelítő megoldások generálásának leállítását, azaz az iterációk leállítását, kilépési feltételhez, feltételekhez kötjük.

Definíció. Az \mathbf{A} mátrixot $\mathbf{P} - \mathbf{Q}$ alakban kifejezve, az \mathbf{A} mátrix egy felbontásának nevezzük.

Legyen az $\mathbf{A} = \mathbf{P} - \mathbf{Q}$ olyan felbontás, ahol \mathbf{P} nemszinguláris, és a $\mathbf{P}\mathbf{y} = \mathbf{c}$ egyenletrendszer egyszerűen megoldható. Írjuk a (4.1) egyenletrendszert

$$\mathbf{P}\mathbf{x} = \mathbf{b} + \mathbf{Q}\mathbf{x} \quad (4.22)$$

alakban. Az \mathbf{x} vektort úgy próbáljuk meghatározni, hogy egy tetszőleges \mathbf{x}_0 becslésből kiindulva egy $\{\mathbf{x}_i\}$ sorozatot generálunk, a

$$\mathbf{P}\mathbf{x}_{i+1} = \mathbf{b} + \mathbf{Q}\mathbf{x}_i \quad (4.23)$$

egyenletrendszer sorozatos megoldásával. [3]

4.2. Tétel. A (4.23) egyenlőségből képzett $\{\mathbf{x}_i\}$ sorozat a (4.1) egyenletrendszer megoldásához konvergál, ha $\|\mathbf{P}^{-1}\mathbf{Q}\| < 1$, tetszőleges normával. [3, 68. oldal, Theorem 4.3]

Bizonyítás. Legyen a hibavektor $\mathbf{e}_i = \mathbf{x}_i - \mathbf{x}$, ahol $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$, a (4.1) egyenletrendszer megoldása. A (4.22) összefüggést kivonva a (4.23) összefüggésből és balról \mathbf{P}^{-1} -el szorozva kapjuk, hogy $\mathbf{e}_{i+1} = \mathbf{P}^{-1}\mathbf{Q}\mathbf{e}_i$.

Legyen $\|P^{-1}Q\| = \alpha$, ekkor $\|e_{i+1}\| \leq \alpha \|e_i\|$, és ebből következik, hogy $\|e_i\| \leq \alpha^i \|e_0\|$. Tehát ha $\alpha < 1$, elég nagy i -re $\|e_i\|$ tetszőlegesen kicsi. Azaz $\lim_{i \rightarrow \infty} x_i = x$. \square

A fentiekben feltételeztük, hogy az A mátrix felbontása, a P és a Q mátrixok i -től függetlenek.

Definíció. Az iterációs eljárást stacionáriusnak nevezzük, ha az együtthatómátrix felbontása nem függ i -től.

Legyen A nonszinguláris mátrix és legyen $A = D - L - U$, ahol D diagonális mátrix, L alsó háromszögmátrix és U felső háromszögmátrix. Néhány példa stacionárius iterációs eljárásokra:

- Jacobi-iteráció: $P = D$, $Q = L + U$,
- Gauss-Seidel eljárás: $P = D - L$, $Q = U$,
- felső relaxációs eljárás: $P = \omega^{-1}D - L$, $Q = U + (\omega^{-1} - 1)D$.

Az ω skalárral a konvergenciát gyorsítjuk. Ez általában azt jelenti, hogy $1 \leq \omega < 2$. Ha $\omega = 1$, a módszer pontosan a Gauss-Seidel eljárás. [3]

A következő fejezetekben két iterációs eljárással fogunk részletesebben foglalkozni, Broyden kutatásainak eredményeire támaszkodva. A felső relaxációs eljárással és a blokk konjugált gradiens módszerrel.

4.3. Sajátértékek és sajátvektorok

Ha az A mátrix négyzetes és nonszinguláris, akkor az $Ax = b$ egyenletrendszernek egyértelmű a megoldása. Ilyen egyenletrendszerek sokszor előfordulnak, amikor valamilyen rendszer statikus viselkedését vizsgáljuk, ahol a rendszer válasza a valamilyen kísérletre a b vektor. Azonban, ha az ilyen rendszerek dinamikus viselkedését akarjuk vizsgálni, akkor meg kell határoznunk azokat a λ skalárokat, amelyekre az $(A - \lambda I)$ mátrix szinguláris. Ezek a λ skalárok a vizsgált rendszer belső tulajdonságaival függenek össze. [3]

Definíció. Legyen A négyzetes mátrix. A λ skalár, ami mellett az $A - \lambda I$ mátrix szinguláris, az A mátrix sajátértéke.

Definíció. Legyen A négyzetes mátrix. Az $x \neq 0$ vektor, ami mellett

$$(A - \lambda I)x = 0, \tag{4.24}$$

az \mathbf{A} mátrix jobboldali sajátvektora. Legyen \mathbf{B} négyzetes mátrix. Az $\mathbf{y} \neq \mathbf{0}$ vektor, ami mellett

$$\mathbf{y}^T(\mathbf{B} - \lambda \mathbf{I}) = \mathbf{0}^T, \quad (4.25)$$

a \mathbf{B} mátrix baloldali sajátvektora.

A (4.24) egyenlőségből

$$\mathbf{A}\mathbf{x} = \lambda\mathbf{x}. \quad (4.26)$$

Definíció. Legyen \mathbf{A} n -ed rendű négyzetes mátrix és $\mathbf{x} \neq \mathbf{0}$ tetszőleges n -ed rendű vektor. Ekkor az

$$\mathbf{x}, \mathbf{A}\mathbf{x}, \mathbf{A}^2\mathbf{x}, \dots$$

sorozatot Krylov sorozatnak nevezzük.

4.3. Tétel. Minden négyzetes mátrixnak van legalább egy sajátértéke. [3, 76. oldal, Theorem 5.1]

Bizonyítás. Legyen \mathbf{A} n -ed rendű négyzetes mátrix és $\mathbf{x} \neq \mathbf{0}$ tetszőleges n -ed rendű vektor. Az $\mathbf{x}, \mathbf{A}\mathbf{x}, \mathbf{A}^2\mathbf{x}, \dots, \mathbf{A}^k\mathbf{x}$ Krylov sorozat vektorai $k > n$ esetén lineárisan összefüggők a 3.1-es tétel miatt, de nem zárhatjuk ki, hogy a sorozat ennél kevesebb vektora is lineárisan összefüggő. Tegyük fel, hogy a sorozat első $r \leq n$ vektora lineárisan független, de az első $r + 1$ darab vektora lineárisan összefüggő. Ekkor léteznek olyan α_i , $i = 0, 1, \dots, r$, nem mind nulla, $r + 1$ darab skalárok, hogy

$$\mathbf{x}\alpha_0 + \mathbf{A}\mathbf{x}\alpha_1 + \dots + \mathbf{A}^r\mathbf{x}\alpha_r = \mathbf{0}. \quad (4.27)$$

Az $\alpha_r \neq 0$, ugyanis ebből a sorozat első $r + 1$ darab vektorának lineáris függetlensége következne, ami ellentmond a hipotézisnek. Ezért a (4.27) egyenlőség átírható a

$$\left(\frac{\alpha_0}{\alpha_r} \mathbf{I} + \frac{\alpha_1}{\alpha_r} \mathbf{A} + \frac{\alpha_2}{\alpha_r} \mathbf{A}^2 + \dots + \mathbf{A}^r \right) \mathbf{x} = \mathbf{0} \quad (4.28)$$

alakra. Egyszerűbben

$$p(\mathbf{A})\mathbf{x} = \mathbf{0},$$

ahol

$$p(\mathbf{A}) = \frac{\alpha_0}{\alpha_r} \mathbf{I} + \frac{\alpha_1}{\alpha_r} \mathbf{A} + \dots + \mathbf{A}^r. \quad (4.29)$$

Legyen

$$p(\xi) = \frac{\alpha_0}{\alpha_r} + \frac{\alpha_1}{\alpha_r} \xi + \dots + \xi^r \quad (4.30)$$

polinom, ahol ξ skalár. Az algebra alaptételéből tudjuk, hogy

$$p(\xi) = (\xi - \lambda_1)(\xi - \lambda_2) \dots (\xi - \lambda_r), \quad (4.31)$$

ahol $\lambda_1, \lambda_1, \dots, \lambda_r$ a (nem feltétlen különböző) gyökei a $p(\xi) = 0$ egyenletnek.

$$p(\mathbf{A}) = (\mathbf{A} - \lambda_1 \mathbf{I})(\mathbf{A} - \lambda_2 \mathbf{I}) \dots (\mathbf{A} - \lambda_r \mathbf{I}). \quad (4.32)$$

Ezt onnan láthatjuk, ha a (4.31) és (4.32) egyenlőségek jobb oldalait kifejtjük, és összehasonlítjuk a ξ^i és \mathbf{A}^i együtthatókat, $0 \leq i \leq r$.

Most megmutatjuk, hogy $(\mathbf{A} - \lambda_1 \mathbf{I})$ szinguláris. Tegyük fel az ellenkezőjét. Ekkor a (4.28) egyenlőséget balról $(\mathbf{A} - \lambda_1 \mathbf{I})^{-1}$ -el megszorozva kapjuk, a (4.29) és (4.32) egyenlőséget felhasználva, hogy

$$(\mathbf{A} - \lambda_2 \mathbf{I})(\mathbf{A} - \lambda_3 \mathbf{I}) \dots (\mathbf{A} - \lambda_r \mathbf{I}) \mathbf{x} = \mathbf{0}, \quad (4.33)$$

ahol most \mathbf{x} $(r-1)$ -ed rendű. A (4.33) egyenlőség úgy is írható, hogy

$$(\beta_0 \mathbf{I} + \beta_1 \mathbf{A} + \dots + \mathbf{A}^{r-1}) \mathbf{x} = \mathbf{0},$$

vagy

$$\mathbf{x} \beta_0 + \mathbf{A} \mathbf{x} \beta_1 + \dots + \mathbf{A}^{r-1} \mathbf{x} = \mathbf{0},$$

ahol $0 \leq i \leq r-2$, és β_i -k a megfelelő együtthatók. Ebből az következik, hogy az $\mathbf{x}, \mathbf{A} \mathbf{x}, \dots$ sorozat első r vektora lineárisan összefüggő, ami ellentmond a hipotézisnek. Ez az ellentmondás garantálja az $(\mathbf{A} - \lambda_1 \mathbf{I})$ mátrix szingularitását, és az \mathbf{A} mátrix legalább egy sajátértékének létezését. \square

4.4. Tétel. *Legyen \mathbf{A} komplex négyzetes mátrix. Ha λ sajátértéke és \mathbf{x} sajátvektora az \mathbf{A} mátrixnak, akkor $\bar{\lambda}$ is sajátértéke és $\bar{\mathbf{x}}$ is sajátvektora \mathbf{A} mátrixnak. [3, 77. oldal, Theorem 5.2]*

Bizonyítás. Mivel $\mathbf{A} \mathbf{z} = \mathbf{z} \lambda$, a komplex konjugáltakat véve $\overline{(\mathbf{A} \mathbf{z})} = \overline{(\mathbf{z} \lambda)}$. $\overline{(\mathbf{A} \mathbf{z})} = \overline{\mathbf{A} \mathbf{z}}$ és $\overline{(\mathbf{z} \lambda)} = \bar{\mathbf{z}} \bar{\lambda}$. \square

Következmény. *Ha \mathbf{A} valós mátrix, és ha van komplex sajátértéke és sajátvektora, ezek mindig komplex konjugált párokban jelentkeznek.*

Bizonyítás. Ha \mathbf{A} valós mátrix, akkor $\mathbf{A} = \bar{\mathbf{A}}$. Az eredmény a tételből közvetlenül következik. \square

4.5. Tétel. *Legyen az \mathbf{A} mátrix pozitív definit. Ekkor az \mathbf{A} mátrix minden sajátértéke pozitív.*

Bizonyítás. Tegyük fel, hogy λ sajátértéke az \mathbf{A} mátrixnak, és \mathbf{A} pozitív definit. Ha $\lambda = 0$, akkor létezik az \mathbf{A} mátrixnak olyan \mathbf{x} sajátvektora, hogy $\mathbf{A} \mathbf{x} = \mathbf{0}$. De ekkor $\mathbf{x}^T \mathbf{A} \mathbf{x} = 0$, ami ellentmondás. Ha $\lambda < 0$, akkor létezik az \mathbf{A} mátrixnak olyan \mathbf{x} sajátvektora, hogy $\mathbf{A} \mathbf{x} = \lambda \mathbf{x}$. De ekkor $\mathbf{x}^T \mathbf{A} \mathbf{x} = \lambda |\mathbf{x}|^2$, ami negatív, így ellentmondás. \square

Bevezetjük a valós ortogonális mátrix komplex megfelelőjét, az unitér mátrixot.

Definíció. Az \mathbf{A} komplex mátrix unitér mátrix, ha $\mathbf{A}^H \mathbf{A} = \mathbf{I}$.

A Hermite-féle mátrixot ezért a valós szimmetrikus mátrix komplex megfelelőjének tekintjük.

Definíció. Legyen \mathbf{A} és \mathbf{P} n -ed rendű mátrix, és legyen \mathbf{P} nonszinguláris. A \mathbf{PAP}^{-1} transzformációt hasonlósági transzformációnak, az \mathbf{A} és \mathbf{PAP}^{-1} mátrixokat hasonlónak nevezzük. Ezenfelül, ha \mathbf{P} mátrix ortogonális (unitér), a \mathbf{PAP}^{-1} hasonlósági transzformációt ortogonális (unitér) transzformációnak nevezzük.

4.6. Tétel. *A hasonlósági transzformáció nem változtatja meg a mátrix sajátértékeit. [3, 79. oldal, Theorem 5.3]*

Bizonyítás. Legyen \mathbf{A} és \mathbf{P} n -ed rendű mátrix, és legyen \mathbf{P} nonszinguláris. Legyen $\mathbf{Ax} = \mathbf{x}\lambda$. Ekkor $\mathbf{PAP}^{-1}\mathbf{Px} = \mathbf{Px}\lambda$. Ezért, ha λ az \mathbf{A} mátrixhoz tartozó sajátérték, és \mathbf{x} a hozzá tartozó sajátvektor, akkor λ a \mathbf{PAP}^{-1} mátrixnak is sajátértéke, és \mathbf{Px} a hozzá tartozó sajátvektora. \square

4.7. Tétel. *Minden k -ad rendű \mathbf{A} mátrixhoz létezik olyan \mathbf{X} unitér mátrix, hogy $\mathbf{X}^H \mathbf{AX}$ első oszlopa a k -ad rendű egységmátrix első oszlopának többszöröse. [3, 79. oldal, Lemma 5.2]*

Bizonyítás. Legyen $\mathbf{Ax} = \lambda\mathbf{x}$, ahol \mathbf{x} normalizált, azaz $\mathbf{x}^H \mathbf{x} = 1$. Legyen \mathbf{X} olyan unitér mátrix, aminek az első oszlopa \mathbf{x} . Ekkor $\mathbf{x} = \mathbf{X}\mathbf{e}_1$, így $\mathbf{AXe}_1 = \mathbf{X}\mathbf{e}_1\lambda$. Balról szorozva \mathbf{X}^H mátrixszal kapjuk, hogy

$$\mathbf{X}^H \mathbf{AXe}_1 = \mathbf{e}_1\lambda. \quad (4.34)$$

\square

Megmutatjuk, hogy minden n -ed rendű négyzetes mátrixot hasonlósági transzformációval felső háromszögmátrixszá alakíthatunk. A felső háromszögmátrix sajátértékei pedig a főátló elemei.

4.8. Tétel. *Minden n -ed rendű \mathbf{A} mátrixhoz létezik egy olyan \mathbf{Q} unitér mátrix, hogy*

$$\mathbf{Q}^H \mathbf{AQ} = \mathbf{U}, \quad (4.35)$$

ahol \mathbf{U} felső háromszögmátrix. [3, 80. oldal, Theorem 5.4 (Schur's theorem)]

Bizonyítás. Indukcióval bizonyítunk. Tegyük fel, hogy létezik olyan \mathbf{Q}_r unitér mátrix, hogy

$$\mathbf{Q}_r^H \mathbf{AQ}_r = \begin{bmatrix} \mathbf{U}_r & \mathbf{B}_r \\ \mathbf{0} & \mathbf{C}_r \end{bmatrix}, \quad (4.36)$$

ahol \mathbf{U}_r r -ed rendű felső háromszögmátrix. Megmutatjuk, hogy ekkor az egyenlőség $(r+1)$ -re is igaz. Mivel \mathbf{C}_r négyzetes, a 4.7-es tétel miatt létezik olyan \mathbf{X}_r unitér mátrix, hogy

$$\mathbf{X}_r^H \mathbf{C}_r \mathbf{X}_r \mathbf{e}_1 = \mathbf{e}_1 \lambda, \quad (4.37)$$

ahol λ a \mathbf{C}_r részmátrix sajátértéke. Legyen

$$\mathbf{P}_r = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{X}_r \end{bmatrix}. \quad (4.38)$$

Mivel \mathbf{P}_r unitér mátrix, a (4.36) és (4.38) egyenlőségekből következik, hogy

$$\mathbf{P}_r \mathbf{Q}_r^H \mathbf{A} \mathbf{Q}_r \mathbf{P}_r = \begin{bmatrix} \mathbf{U}_r & \mathbf{B}_r \mathbf{X}_r \\ \mathbf{0} & \mathbf{X}_r \mathbf{C}_r \mathbf{X}_r \end{bmatrix}. \quad (4.39)$$

Így a (4.37) egyenlőség miatt a (4.39) egyenlőség jobb oldala írható úgy, hogy

$$\begin{bmatrix} \mathbf{U}_{r+1} & \mathbf{B}_{r+1} \\ \mathbf{0} & \mathbf{C}_{r+1} \end{bmatrix}.$$

Legyen $\mathbf{Q}_{r+1} = \mathbf{Q}_r \mathbf{P}_r$, ekkor \mathbf{Q}_{r+1} unitér, mert két unitér mátrix szorzata unitér. Így a (4.36) egyenlőség $(r+1)$ -re is igaz. Mivel a (4.36) egyenlőség $r=1$ esetén igaz a 4.7-es tétel miatt, így minden r -re, $1 \leq r \leq n-1$, igaz. Ha $r=n-1$ a (4.36) egyenlőség jobb oldala felső háromszög mátrix, ezt \mathbf{U} -val jelölve és \mathbf{Q}_{n-1} -et \mathbf{Q} -val jelölve, kapjuk a tétel állítását. \square

A tétel következményeit bizonyítás nélkül közöljük, a bizonyítások megtalálhatók [3] 80-82. oldalán.

Következmény. Az n -ed rendű \mathbf{A} négyzetes mátrixnak n darab sajátértéke van.

A tétel következő következménye előtt bevezetjük a normálmátrix fogalmát.

Definíció. Az \mathbf{A} mátrix normálmátrix, ha $\mathbf{A}^H \mathbf{A} = \mathbf{A} \mathbf{A}^H$.

Következmény. Az n -ed rendű \mathbf{A} négyzetes mátrix diagonális mátrixszá transzformálható egy unitér transzformációval akkor és csak akkor, ha az \mathbf{A} mátrix normálmátrix.

Következmény. Ha az \mathbf{A} Hermite-féle mátrix, akkor az \mathbf{U} mátrix a (4.35) egyenlőségben valós diagonális mátrix, és ezért az \mathbf{A} mátrix sajátértékei valósak.

Következmény. Ha az \mathbf{A} mátrix valós, és minden sajátértéke valós, akkor az \mathbf{A} mátrix valós felső háromszögmátrixszá alakítható egy ortogonális transzformációval.

Következmény. Ha az \mathbf{A} mátrix valós és szimmetrikus, akkor az \mathbf{A} mátrix valós diagonális mátrixszá alakítható egy ortogonális transzformációval.

Megmutattuk, hogy az n -ed rendű \mathbf{A} négyzetes mátrixnak n darab sajátértéke van, de ezek a sajátértékek nem feltétlen különbözőek. Most megvizsgáljuk hány darab különböző, azaz lineárisan független sajátvektora van egy mátrixnak. Kiemelt fontosságú, hogy egy n -ed rendű négyzetes mátrixnak van-e n darab lineárisan független sajátvektora. Ha igen, akkor ezek a lineárisan független sajátvektorok bázist alkotnak.

Definíció. Ha egy mátrix sajátvektorai nem alkotnak bázist, akkor a mátrixot defektív mátrixnak nevezzük.

Definíció. A mátrixot amelynek oszlopai egy nemdefektív mátrix sajátvektorai, modálmátrixnak nevezzük.

A modálmátrix nem egyértelmű. Egyrészt nincs meghatározva a mátrix oszlopainak sorrendje, másrészt a sajátvektorok tetszőlegesen méretezhetők.

4.9. Tétel. Legyen \mathbf{A} n -ed rendű mátrix és legyen m darab különböző sajátértéke, $m \leq n$. Az \mathbf{A} mátrixnak legalább m darab lineárisan független sajátvektora van. [3, 86. oldal, Theorem 5.6]

Bizonyítás. A tételt először az $\mathbf{U} = [u_{ij}]$ felső háromszögmátrixra bizonyítjuk. Jelölje λ_i az \mathbf{U} mátrix sajátértékeit, λ_k pedig a többszörös sajátértékeit. Legyen

$$\mathbf{U} = \begin{bmatrix} \mathbf{U}_{11} & \mathbf{U}_{12} \\ \mathbf{0} & \mathbf{U}_{22} \end{bmatrix}, \quad (4.40)$$

ahol \mathbf{U}_{11} és \mathbf{U}_{22} négyzetes mátrixok. Particionáljuk úgy az \mathbf{U} mátrixot, hogy \mathbf{U}_{11} mátrixnak ne legyen diagonális eleme λ_k , az \mathbf{U}_{22} mátrixnak pedig legyen az első diagonális eleme λ_k . Tekintsük az $(\mathbf{U} - \lambda_k \mathbf{I}) = \mathbf{0}$ egyenlőséget, amit írhatunk az

$$\begin{bmatrix} \mathbf{U}_{11} - \lambda_k \mathbf{I} & \mathbf{U}_{12} \\ \mathbf{0} & \mathbf{U}_{22} - \lambda_k \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \quad (4.41)$$

vagy a

$$(\mathbf{U}_{11} - \lambda_k \mathbf{I})\mathbf{x}_1 = -\mathbf{U}_{12}\mathbf{x}_2 \quad (4.42)$$

és

$$(\mathbf{U}_{22} - \lambda_k \mathbf{I})\mathbf{x}_2 = \mathbf{0} \quad (4.43)$$

alakban. A particionálás miatt az $(\mathbf{U}_{22} - \lambda_k \mathbf{I})$ mátrix első oszlopa nulla, ezért $\mathbf{x}_2 = \mathbf{e}_1$, ahol \mathbf{e}_1 a megfelelő rendű egységmátrix első oszlopa. A particionálás miatt az

$(\mathbf{U}_{11} - \lambda_k \mathbf{I})\mathbf{x}_1$ nonszinguláris. Ekkor az (4.42) egyenlőségből

$$\mathbf{x}_1 = -(\mathbf{U}_{11} - \lambda_k \mathbf{I})^{-1} \mathbf{U}_{12} \mathbf{e}_2. \quad (4.44)$$

A (4.41) egyenlőségből látszik, hogy ha \mathbf{U}_{11} $(r-1)$ -ed rendű, akkor az \mathbf{x} sajátvektor az r -edik oszlopa valamilyen n -ed rendű felső egység háromszögmátrixnak. Így az \mathbf{U} mátrix m darab különböző sajátértékéhez tartozik különböző sajátvektor, amelyek egy felső egység háromszögmátrix különböző oszlopai. Mivel a felső egység háromszögmátrix nonszinguláris, a sajátvektorok lineárisan függetlenek.

Hogy a tételt az általános \mathbf{A} mátrixra belássuk, legyen

$$\mathbf{Q}^H \mathbf{A} \mathbf{Q} = \mathbf{U}, \quad (4.45)$$

ahol \mathbf{Q} unitér mátrix. Legyen \mathbf{V} $n \times m$ -es mátrix a $\lambda_1, \lambda_2, \dots, \lambda_m$ különböző sajátértékekhez tartozó lineárisan független sajátvektorok mátrixa. Legyen $\mathbf{\Lambda}_1 = \text{diag}(\lambda_i)$ az m -ed rendű diagonális mátrix a $\lambda_1, \lambda_2, \dots, \lambda_m$ különböző sajátértékekkel. Ekkor

$$\mathbf{U} \mathbf{V} = \mathbf{V} \mathbf{\Lambda}_1. \quad (4.46)$$

A (4.45) és a (4.46) egyenlőségekből következik, hogy

$$\mathbf{A} \mathbf{Q} \mathbf{V} = \mathbf{Q} \mathbf{V} \mathbf{\Lambda}_1.$$

Azaz a $\mathbf{Q} \mathbf{V}$ mátrix oszlopai az \mathbf{A} mátrixnak a $\lambda_1, \lambda_2, \dots, \lambda_m$ különböző sajátértékeihez tartozó lineárisan független sajátvektorok mátrixa. Mivel \mathbf{V} mátrix oszlopai lineárisan függetlenek és \mathbf{Q} unitér mátrix, így a $\mathbf{Q} \mathbf{V}$ mátrix nonszinguláris, tehát a sajátvektorok lineárisan függetlenek. \square

A tétel következményeit bizonyítás nélkül közöljük, a bizonyítások megtalálhatók [3] 88. oldalán.

Következmény. *Ha egy n -ed rendű mátrixnak n darab különböző sajátértéke van, akkor nem defektív mátrix.*

Következmény. *Ha egy mátrix sajátértéke egyszeres, a hozzá tartozó sajátvektor egyértelmű, a méretezést leszámítva.*

Következmény. *Ha egy mátrix sajátértéke r -szeres, akkor legfeljebb r darab lineárisan független sajátvektor tartozik hozzá.*

A Iteratív eljárások fejezetben láttuk, hogy az iteratív eljárások konvergenciája összefügg a mátrixok normájával. Megvizsgálunk néhány, a sajátérték, a spektrális sugár és a mátrix norma közötti összefüggést.

Definíció. Legyen \mathbf{A} n -ed rendű mátrix, λ_i sajátvektorokkal, $1 \leq i \leq n$. Az \mathbf{A} mátrix spektrális sugara $\rho(\mathbf{A}) = \max_i |\lambda_i|$.

4.10. Tétel. Legyen \mathbf{A} négyzetes mátrix. Ekkor $\rho(\mathbf{A}) \leq \|\mathbf{A}\|$ tetszőleges normával. [3, 89. oldal, Theorem 5.7]

Bizonyítás. Legyen \mathbf{x} az \mathbf{A} mátrix egy sajátvektora és legyen λ a hozzá tartozó sajátérték, azaz $\mathbf{Ax} = \lambda\mathbf{x}$. Így a normák tulajdonságai alapján

$$\|\lambda\mathbf{x}\| = \|\lambda\| \|\mathbf{x}\| = \|\mathbf{Ax}\| \leq \|\mathbf{A}\| \|\mathbf{x}\|.$$

Mivel $\mathbf{x} \neq \mathbf{0}$, így $|\lambda| \leq \|\mathbf{A}\|$. □

4.11. Tétel. Legyen \mathbf{A} $m \times n$ -es valós mátrix. Ekkor

$$\|\mathbf{A}\|_2 = [\rho(\mathbf{A}^T \mathbf{A})]^{1/2}.$$

[3, 89. oldal, Theorem 5.8]

Bizonyítás. Az indukált mátrixnorma definíciójából,

$$\|\mathbf{A}\|_2^2 = \left(\max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{Ax}\|_2}{\|\mathbf{x}\|_2} \right)^2 = \max_{\mathbf{x} \neq \mathbf{0}} \left(\frac{\|\mathbf{Ax}\|_2}{\|\mathbf{x}\|_2} \right)^2 = \max_{\mathbf{x} \neq \mathbf{0}} \left(\frac{\mathbf{x}^T \mathbf{A}^T \mathbf{Ax}}{\mathbf{x}^T \mathbf{x}} \right). \quad (4.47)$$

Mivel az $\mathbf{A}^T \mathbf{A}$ mátrix valós és szimmetrikus, a sajátvektorai valósak és az \mathbf{X} modál-mátrixa ortogonális, a 4.8-as tétel következményei miatt. Legyen $\mathbf{\Lambda}$ a $\mathbf{A}^T \mathbf{A}$ mátrix sajátvektorainak diagonális mátrixa, azaz $\mathbf{A}^T \mathbf{Ax} = \mathbf{X}\mathbf{\Lambda}$. Legyen $\mathbf{x} = \mathbf{Xz}$. Mivel $\mathbf{X}^T \mathbf{X} = \mathbf{I}$, következik, hogy $\mathbf{x}^T \mathbf{A}^T \mathbf{Ax} = \mathbf{z}^T \mathbf{\Lambda} \mathbf{z}$, és $\mathbf{x}^T \mathbf{x} = \mathbf{z}^T \mathbf{z}$. Így a (4.47) egyenlőség alakja

$$\|\mathbf{A}\|_2^2 = \max_{\mathbf{z} \neq \mathbf{0}} \left(\frac{\mathbf{z}^T \mathbf{\Lambda} \mathbf{z}}{\mathbf{z}^T \mathbf{z}} \right) = \max_{\mathbf{z} \neq \mathbf{0}} \left[\left(\sum_{i=1}^n \lambda_i z_i^2 \right) / \left(\sum_{i=1}^n z_i^2 \right) \right]. \quad (4.48)$$

Legyen

$$\theta_i = \frac{z_i^2}{\sum_{j=1}^n z_j^2},$$

ezt a (4.48) egyenlőségbe helyettesítve kapjuk, hogy

$$\|\mathbf{A}\|_2^2 = \max_i \sum_{i=1}^n \theta_i \lambda_i,$$

ahol $\theta \geq 0$ és $\sum_{i=1}^n \theta_i = 1$. A $\sum_i \theta_i \lambda_i$ kifejezést maximalizáljuk. A legnagyobb pozitív λ_i -hoz tartozó θ_i -t válasszuk egynek, az összes többi θ_i pedig nullának. Mivel az

$\mathbf{A}^T \mathbf{A}$ mátrixhoz tartozó összes sajátérték pozitív, a legnagyobb pozitív sajátértéke az $\mathbf{A}^T \mathbf{A}$ mátrixnak a spektrális sugara. Azaz $\|\mathbf{A}\|_2^2 = \rho(\mathbf{A}^T \mathbf{A})$. \square

Következmény. Ha $\mathbf{A} = \mathbf{A}^T$, $\|\mathbf{A}\|_2 = \rho(\mathbf{A})$.

Bizonyítás. Ha $\mathbf{A} = \mathbf{A}^T$, akkor $\mathbf{A}^T \mathbf{A} = \mathbf{A}^2$. Ha λ sajátértéke az \mathbf{A} mátrixnak, akkor λ^2 sajátértéke az \mathbf{A}^2 mátrixnak. \square

Következmény. $\|\mathbf{A}\|_2^2 \leq \|\mathbf{A}\|_1 \|\mathbf{A}\|_\infty$.

Bizonyítás. Legyen \mathbf{x} az $\mathbf{A}^T \mathbf{A}$ mátrix $\rho(\mathbf{A}^T \mathbf{A})$ sajátértékéhez tartozó sajátvektora. Ekkor

$$\mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{x} \rho(\mathbf{A}^T \mathbf{A}),$$

és a tételt alkalmazva

$$\|\mathbf{x} \rho(\mathbf{A}^T \mathbf{A})\|_\infty = \|\mathbf{x}\|_\infty \|\mathbf{A}\|_2^2 = \|\mathbf{A}^T \mathbf{A} \mathbf{x}\|_\infty \leq \|\mathbf{A}^T\|_\infty \|\mathbf{A}\|_\infty \|\mathbf{x}\|_\infty,$$

ahol $\|\mathbf{x}\|_\infty \neq 0$ és $\|\mathbf{A}^T\|_\infty = \|\mathbf{A}\|_1$. \square

Definíció. Legyen $\mathbf{A} = [a_{ij}]$ n -ed rendű mátrix. Ekkor az \mathbf{A} mátrix nyoma $tr(\mathbf{A}) = \sum_{i=1}^n a_{ii}$.

4.12. Tétel. *A mátrixhoz tartozó sajátértékek összege egyenlő a mátrix nyomával. [3, 91. oldal, Theorem 5.9]*

Bizonyítás. Bármilyen \mathbf{B} és \mathbf{C} mátrixokra, amelyeknek a \mathbf{BC} és \mathbf{CB} szorzatuk definiálva van, $tr(\mathbf{BC}) = tr(\mathbf{CB})$. A 4.8-as tételből tudjuk, hogy $\mathbf{Q}^H \mathbf{A} \mathbf{Q} = \mathbf{U}$. Ezért

$$tr(\mathbf{U}) = tr((\mathbf{Q}^H \mathbf{A} \mathbf{Q}) \mathbf{Q}) = tr(\mathbf{Q}(\mathbf{Q}^H \mathbf{A})) = tr(\mathbf{A}),$$

az \mathbf{U} mátrix nyoma pedig megegyezik az \mathbf{U} mátrix sajátvektorainak összegével. \square

4.4. Kovergens mátrixok

A különböző iteratív eljárások konvergenciája gyakran egy \mathbf{A}^r mátrix viselkedésétől függ, miközben r tart a végtelenhez. Bizonyos iteratív eljárások akkor és csak akkor konvergensek, ha \mathbf{A}^r a nullmátrixhoz konvergál, ahogy növekszik r . [3]

Definíció. Az n -ed rendű \mathbf{A} mátrix konvergens, ha $\lim_{r \rightarrow \infty} \|\mathbf{A}^r\| = 0$.

4.13. Tétel. *Legyen \mathbf{U} felső háromszögmátrix és $\rho(\mathbf{U}) < 1$. Ekkor létezik \mathbf{B} diagonális mátrix, amelyre $\|\mathbf{B} \mathbf{U} \mathbf{B}^{-1}\|_\infty < 1$ és $\|\mathbf{B} \mathbf{U} \mathbf{B}^{-1}\|_1 < 1$. [3, 92. oldal, Theorem 5.10]*

Bizonyítás. Legyen $\mathbf{U} = [u_{ij}]$ n -ed rendű felső háromszögmátrix. Legyen $\mathbf{B} = \text{diag}(1, \beta, \beta^2, \dots, \beta^{n-1})$, ahol válasszunk $\beta > 1$ skalárt úgy, hogy

$$\beta^{j-i} > \frac{(n-1)|u_{ij}|}{1-\rho(\mathbf{U})}, \quad j > i \quad (4.49)$$

teljesüljön. Legyen $\mathbf{V} = [v_{ij}] = \mathbf{B}\mathbf{U}\mathbf{B}^{-1}$. Ekkor $j > i$ -re $v_{ij} = u_{ij}/\beta^{j-i}$ és a (4.49) egyenlőtlenség miatt

$$|v_{ij}| < \frac{1-\rho(\mathbf{U})}{n-1}, \quad j > i. \quad (4.50)$$

De $|v_{ii}| = |u_{ii}| < \rho(\mathbf{U})$ és $v_{ij} = 0$ minden $i > j$ -re. Ezért

$$\|\mathbf{V}\|_{\infty} \leq \rho(\mathbf{U}) + \max_i \sum_{j=i+1}^n |v_{ij}|,$$

és így a (4.50) egyenlőtlenség miatt

$$\|\mathbf{V}\|_{\infty} \leq \rho(\mathbf{U}) + (n-1) \left[\frac{1-\rho(\mathbf{U})}{n-1} \right] = 1.$$

Hasonlóan, $\|\mathbf{V}\|_1 < 1$. □

Következmény. A tétel feltételei mellett $\|\mathbf{B}\mathbf{U}\mathbf{B}^{-1}\|_2 < 1$.

Bizonyítás. A 4.11-es tételből következik. □

A tétel fontossága, hogy ha az \mathbf{A} mátrix spektrális sugara kisebb mint 1, akkor mindig létezik olyan norma, amiben az \mathbf{A} mátrix normája kisebb mint 1.

4.14. Tétel. Az n -ed rendű \mathbf{A} mátrix konvergenciájának szükséges és elégséges feltétele, hogy $\rho(\mathbf{A}) < 1$. [3, 93. oldal, Theorem 5.11]

Bizonyítás. A tételt először egy $\mathbf{U} = [u_{ij}]$ felső háromszögmátrixra bizonyítjuk, majd a 4.8-as tétellel kiterjesztjük a bizonyítást tetszőleges négyzetes mátrixra. Az i -edik diagonális eleme \mathbf{U}^r mátrixnak u_{ii}^r , ezért \mathbf{U} nem lehet konvergens, ha $u_{ii} \geq 1$ bármilyen i -re. Ekkor ugyanis az \mathbf{U} mátrix hatványozásakor az aktuális u_{ii}^r hatványok folyamatosan nőnek. Ha $u_{ii} \leq 1$, az \mathbf{U} mátrix hatványozásakor az aktuális u_{ii}^r hatványok folyamatosan csökkennek. Mivel $\rho(\mathbf{U}) = \max_i |u_{ii}|$, ezért $\rho(\mathbf{U}) < 1$ szükséges feltétele a konvergenciának. Bizonyítjuk a feltétel elégségességét. A 4.13-es és a 4.11-es tétel következménye, hogy ha $\rho(\mathbf{U}) < 1$, akkor létezik olyan \mathbf{B} diagonális mátrix, hogy $\mathbf{V} = \mathbf{B}\mathbf{U}\mathbf{B}^{-1}$ és $\|\mathbf{V}\|_2 < 1$. Ezért $\mathbf{U} = \mathbf{B}^{-1}\mathbf{V}\mathbf{B}$ és $\mathbf{U}^r = \mathbf{B}^{-1}\mathbf{V}^r\mathbf{B}$, mivel $\|\mathbf{B}\|_2 \|\mathbf{B}^{-1}\|_2 = \beta^{n-1}$. Ezekből következik, hogy $\|\mathbf{U}^r\|_2 \leq \beta^{n-1} \|\mathbf{V}\|_2^r$. Mivel $\|\mathbf{V}\|_2 < 1$ és $\lim_{r \rightarrow \infty} \|\mathbf{U}^r\| = 0$, az \mathbf{U} mátrix konvergens.

Általános esetben, a 4.8-as tételből $\mathbf{Q}^H \mathbf{A} \mathbf{Q} = \mathbf{U}$, ezért $\mathbf{A} = \mathbf{Q} \mathbf{U} \mathbf{Q}^H$ és $\mathbf{A}^r = \mathbf{Q} \mathbf{U}^r \mathbf{Q}^H$. Mivel $\|\mathbf{Q}\|_2 \|\mathbf{Q}^H\|_2 = 1$, $\|\mathbf{A}^r\|_2 = \|\mathbf{U}^r\|_2$. Ebből az egyenlőségből következik, hogy a tétel általános esetben is igaz.

□

5. Felső relaxációs eljárás (SOR)

Stacionárius iterációs módszereknél minden olyan módosítás ami csökkenti a spektrális sugarat, meggyorsítja a konvergenciáját a lineáris egyenletrendszer megoldásának. A SOR módszer az egyik tagja annak a nagy módszer családnak, amelyet az iterációs eljárások konvergenciájának gyorsítására dolgoztak ki [14].

A SOR módszer konvergencia kritériumát szimmetrikus együtthatómátrixok esetén Alexander Ostrowski dolgozta ki 1954-ben [12]. Broyden elégséges konvergencia feltételeket adott 1964-ben, szimmetrikus és nonszimmetrikus együtthatómátrixok esetére is [2]. Ebben a fejezetben röviden bemutatjuk a SOR módszert, és összefoglaljuk Broyden a SOR módszer konvergenciájával kapcsolatos eredményeit [2].

5.1. Általános konvergencia eredmények

Legyen

$$\mathbf{M} \mathbf{x} = \mathbf{c}$$

lineáris egyenletrendszer, ahol \mathbf{M} nonszinguláris, valós mátrix, \mathbf{c} valós vektor, \mathbf{x} ismeretlen vektor. Legyen

$$\mathbf{A} = \mathbf{D} \mathbf{M}, \quad \mathbf{b} = \mathbf{D} \mathbf{c}, \quad (5.1)$$

ahol \mathbf{D} diagonális mátrix. Így

$$\mathbf{A} \mathbf{x} = \mathbf{b}. \quad (5.2)$$

A \mathbf{D} mátrixot válasszuk úgy, hogy az \mathbf{A} mátrix főátlójában lévő összes elem egységelem legyen. Ekkor az \mathbf{A} mátrixot kifejezhetjük három mátrix összegeként,

$$\mathbf{A} = \mathbf{I} + \mathbf{L} + \mathbf{U}, \quad (5.3)$$

ahol \mathbf{I} az egységmátrix, \mathbf{L} alsó háromszögmátrix, \mathbf{U} pedig felső háromszögmátrix.

Legyen \mathbf{x}_{i+1} és \mathbf{x}_i egymást követő közelítő megoldása az (5.2) egyenletrendszernek. Ekkor a SOR módszer szerint

$$\mathbf{x}_{i+1} = \mathbf{x}_i - \omega [(\mathbf{I} + \mathbf{U}) \mathbf{x}_i + \mathbf{L} \mathbf{x}_{i+1} - \mathbf{b}], \quad (5.4)$$

ahol ω skalár. Az (5.3) egyenlőséget felhasználva, az (5.4) egyenlőségből az \mathbf{U} mátrix eliminálásával kapjuk, hogy

$$(\mathbf{I} + \omega \mathbf{L})(\mathbf{x}_{i+1} - \mathbf{x}_i) = -\omega(\mathbf{A}\mathbf{x}_i - \mathbf{b}). \quad (5.5)$$

Definíció. Az \mathbf{x}_i közelítő megoldáshoz tartozó maradéktag $\boldsymbol{\epsilon}_i = \mathbf{A}\mathbf{x}_i - \mathbf{b}$.

Ebből következik, hogy

$$\boldsymbol{\epsilon}_{i+1} = \mathbf{A}\mathbf{x}_{i+1} - \mathbf{b},$$

és

$$\mathbf{A}^{-1}(\boldsymbol{\epsilon}_{i+1} - \boldsymbol{\epsilon}_i) = \mathbf{x}_{i+1} - \mathbf{x}_i.$$

Így az (5.5) egyenlőségből

$$(\mathbf{I} + \omega \mathbf{L})\mathbf{A}^{-1}(\boldsymbol{\epsilon}_{i+1} - \boldsymbol{\epsilon}_i) = -\omega\boldsymbol{\epsilon}_i. \quad (5.6)$$

Az (5.6) egyenlőségből

$$\boldsymbol{\epsilon}_{i+1} = [\mathbf{I} - \omega \mathbf{A}(\mathbf{I} + \omega \mathbf{L})^{-1}] \boldsymbol{\epsilon}_i. \quad (5.7)$$

Írjuk ezt az egyenlőséget az

$$\boldsymbol{\epsilon}_{i+1} = (\mathbf{I} - \mathbf{B})\boldsymbol{\epsilon}_i$$

egyszerűbb formában. [2]

5.1. Tétel. *Azoknál az iteratív eljárásoknál, ahol a közelítő megoldások maradéktagja kifejezhető az*

$$\boldsymbol{\epsilon}_{i+1} = (\mathbf{I} - \mathbf{B})\boldsymbol{\epsilon}_i$$

alakban, a konvergencia elégséges feltétele, hogy létezzen olyan \mathbf{S} és \mathbf{G} mátrix, hogy

$$\mathbf{S} > 0$$

és

$$\mathbf{G} = \mathbf{B}^T \mathbf{S} + \mathbf{S} \mathbf{B} - \mathbf{B}^T \mathbf{S} \mathbf{B} > 0.$$

[2, 137. oldal, Theorem 1]

Bizonyítás. Legyen $f_i = \boldsymbol{\epsilon}_i^T \mathbf{S} \boldsymbol{\epsilon}_i$. Mivel a hipotézis szerint $\mathbf{S} > 0$, $f_i > 0$, minden $\boldsymbol{\epsilon}_i \neq 0$ esetén. A konvergencia szükséges és elégséges feltétele, ha $f_i \rightarrow 0$, miközben $i \rightarrow \infty$.

$$\begin{aligned}
f_{i+1} &= \boldsymbol{\epsilon}_{i+1}^T \mathbf{S} \boldsymbol{\epsilon}_{i+1} \\
&= \boldsymbol{\epsilon}_i^T (\mathbf{I} - \mathbf{B}^T) \mathbf{S} (\mathbf{I} - \mathbf{B}) \boldsymbol{\epsilon}_i \\
&= \boldsymbol{\epsilon}_i^T (\mathbf{S} - \mathbf{G}) \boldsymbol{\epsilon}_i.
\end{aligned}$$

Legyen $\phi_i = \boldsymbol{\epsilon}_i^T \mathbf{G} \boldsymbol{\epsilon}_i$, ekkor $f_{i+1} = f_i - \phi_i$. Egy elégséges feltétele annak, hogy $f_i \rightarrow 0$, miközben $i \rightarrow \infty$, ha létezik olyan k konstans, hogy

$$\phi_i \geq k f_i, \quad (5.8)$$

mert ekkor $f_{i+1} \leq (1 - k)f_i$ és az f_i sorozat konvergál.

Az \mathbf{S} mátrix minden sajátértéke valós és pozitív, mert $\mathbf{S} > 0$ (4.5-ös tétel). Ha λ_{\max} az \mathbf{S} mátrix legnagyobb sajátértéke, akkor

$$f_i \leq \lambda_{\max} \boldsymbol{\epsilon}_i^T \boldsymbol{\epsilon}_i. \quad (5.9)$$

Mivel $\mathbf{S} = \mathbf{S}^T$, a \mathbf{G} mátrix szimmetrikus és a sajátértékei valósak. Legyen a \mathbf{G} mátrix legkisebb sajátértéke μ_{\min} . Ekkor $\phi_i \geq \mu_{\min} \boldsymbol{\epsilon}_i^T \boldsymbol{\epsilon}_i$. Az (5.9) egyenlőségből

$$\phi_i \geq \frac{\mu_{\min}}{\lambda_{\max}} f_i.$$

Ha $\mathbf{G} > 0$, akkor $\mu_{\min} > 0$, és teljesül az (5.8) feltétel. \square

Következmény. Ha $\mathbf{S} > 0$ és $\mathbf{G} \leq 0$, azaz ha \mathbf{S} mátrix pozitív definit és a \mathbf{G} mátrix negatív szemidefinit, akkor az iteratív eljárás sohasem konvergál.

Most az 5.1-es tételt alkalmazzuk a SOR módszerre. Az 5.7 egyenlőség miatt, $\mathbf{B} = \omega \mathbf{A}(\mathbf{I} + \omega \mathbf{L})^{-1}$, ezért a SOR módszer konvergenciájának egy elégséges feltétele, hogy létezik olyan $\mathbf{S} > 0$ mátrix, hogy

$$\omega(\mathbf{I} + \omega \mathbf{L}^T)^{-1} \mathbf{A}^T \mathbf{S} + \omega \mathbf{S} \mathbf{A}(\mathbf{I} + \omega \mathbf{L})^{-1} - \omega^2 (\mathbf{I} + \omega \mathbf{L}^T)^{-1} \mathbf{A}^T \mathbf{S} \mathbf{A}(\mathbf{I} + \omega \mathbf{L})^{-1} > 0. \quad (5.10)$$

Ezt a feltételt egyszerűsíthetjük a következő tétel felhasználásával.

5.2. Tétel. A $\mathbf{P} > 0$ szükséges és elégséges feltétele, hogy $\mathbf{Q}^T \mathbf{P} \mathbf{Q} > 0$, ahol \mathbf{Q} tetszőleges nemszinguláris mátrix. [2, 138. oldal, Lemma]

Bizonyítás. Legyen $\mathbf{Q} \mathbf{x} = \mathbf{z}$. Ekkor $\mathbf{x}^T \mathbf{Q}^T \mathbf{P} \mathbf{Q} \mathbf{x} = \mathbf{z}^T \mathbf{P} \mathbf{z}$. Mivel \mathbf{Q} mátrix nemszinguláris, minden nem nulla \mathbf{x} vektorra létezik nem nulla \mathbf{z} vektor, és fordítva. \square

Mivel $(\mathbf{I} + \omega \mathbf{L})$ nemszinguláris, ezért az (5.10) elégséges feltétel a rövidebb

$$\omega [\mathbf{A}^T \mathbf{S}(\mathbf{I} + \omega \mathbf{L}) + (\mathbf{I} + \omega \mathbf{L}^T) \mathbf{S} \mathbf{A} - \omega \mathbf{A}^T \mathbf{S}] > 0 \quad (5.11)$$

formában írható [2].

5.2. Szimmetrikus mátrixok konvergenciája

Tegyük fel, hogy az \mathbf{M} mátrix szimmetrikus. Legyen $\mathbf{S} = \mathbf{D}^{-1}\mathbf{M}^{-1}\mathbf{D}^{-1}$. Ekkor az (5.10) feltétel alakja,

$$\omega [\mathbf{D}^{-1}(\mathbf{I} + \omega\mathbf{L}) + (\mathbf{I} + \omega\mathbf{L}^T)\mathbf{D}^{-1} - \omega\mathbf{M}] > 0. \quad (5.12)$$

Az (5.1) és (5.3) egyenlőségekből, és \mathbf{M} feltételezett szimmetrikussága miatt,

$$\mathbf{M} = \mathbf{D}^{-1}(\mathbf{I} + \mathbf{L} + \mathbf{U}) = (\mathbf{I} + \mathbf{L}^T + \mathbf{U}^T)\mathbf{D}^{-1}. \quad (5.13)$$

Az \mathbf{M} mátrix az (5.13) egyenlőség szerinti két reprezentációjának felső háromszög partícióit egyenlővé téve kapjuk, hogy

$$\mathbf{D}^{-1}\mathbf{U} = \mathbf{L}^T\mathbf{D}^{-1}. \quad (5.14)$$

Így az (5.11) feltétel az

$$\omega(2 - \omega)\mathbf{D}^{-1} > 0 \quad (5.15)$$

alakra egyszerűsödik. Ha $\mathbf{M} > 0$ akkor az 5.2-es tételből következik, hogy $\mathbf{S} > 0$, és mivel $\mathbf{M} > 0$, ezért $\mathbf{D} > 0$, és az (5.15) feltétel teljesül. Ezért a SOR módszer konvergálni fog, ha $\mathbf{M} > 0$ és $0 < \omega < 2$. Azonban ha $\omega < 0$, vagy $\omega > 2$, a \mathbf{G} mátrix negatív definitté válik, és a 5.1-es tétel következménye miatt a SOR módszer nem fog konvergálni. [2]

5.3. Nemszimmetrikus mátrixok konvergenciája

Legyen $\mathbf{S} = (\mathbf{A}^T)^{-1}\mathbf{A}^{-1}$ az (5.11) feltételben. Az \mathbf{A} mátrix nemszinguláris, ezért az $(\mathbf{A}^T)^{-1}\mathbf{A}^{-1}$ mátrix pozitív definit. Ekkor a konvergencia feltétele,

$$\omega\mathbf{A}^{-1}(\mathbf{I} + \omega\mathbf{L}) + \omega(\mathbf{I} + \omega\mathbf{L}^T) + (\mathbf{A}^T)^{-1} - \omega^2\mathbf{I} > 0.$$

Az \mathbf{A} mátrix nemszinguláris, ezért az 5.2-es tételt felhasználva, a feltétel a

$$\omega(\mathbf{I} + \omega\mathbf{L})\mathbf{A}^T + \omega\mathbf{A}(\mathbf{I} + \omega\mathbf{L}^T) - \omega^2\mathbf{A}\mathbf{A}^T > 0$$

alakra egyszerűsödik. Az \mathbf{A} mátrixot felbontva (5.3) szerint és egyszerűsítve,

$$\omega(\mathbf{A} + \mathbf{A}^T) - \omega^2 [(\mathbf{I} + \mathbf{U})(\mathbf{I} + \mathbf{U}^T) - \mathbf{L}\mathbf{L}^T] > 0. \quad (5.16)$$

Legyen $\mathbf{S} = \mathbf{I}$ az (5.11) feltételben. Ekkor a konvergencia elégséges feltétele

$$\omega \mathbf{A}^T (\mathbf{I} + \omega \mathbf{L}) + \omega (\mathbf{I} + \omega \mathbf{L}^T) \mathbf{A} - \omega^2 \mathbf{A}^T \mathbf{A} > 0,$$

az \mathbf{A} mátrixot felbontva (5.3) szerint és egyszerűsítve,

$$\omega (\mathbf{A} + \mathbf{A}^T) - \omega^2 [(\mathbf{I} + \mathbf{U}^T)(\mathbf{I} + \mathbf{U}) - \mathbf{L}^T \mathbf{L}] > 0. \quad (5.17)$$

Megmutatjuk, hogy ha $\mathbf{A} + \mathbf{A}^T > 0$, akkor létezik olyan pozitív ω , ami kielégíti az (5.16) és az (5.17) feltételeket. [2]

5.3. Tétel. *Ha $\mathbf{P} > 0$ és $\mathbf{Q} = \mathbf{Q}^T$, akkor létezik olyan pozitív ω , hogy $\mathbf{P} + \omega \mathbf{Q} > 0$. [2, 139. oldal, Theorem 2]*

Bizonyítás. Legyen $f_1 = \mathbf{x}^T \mathbf{P} \mathbf{x}$. Mivel $\mathbf{P} > 0$, minden sajátértéke valós és pozitív (4.5-ös tétel). Jelölje λ_{\min} a \mathbf{P} mátrix legkisebb sajátértékét. Ekkor

$$f_1 \geq \lambda_{\min} \mathbf{x}^T \mathbf{x}.$$

A \mathbf{Q} mátrix szimmetrikus, ezért minden sajátértéke valós. Jelölje a \mathbf{Q} mátrix legkisebb sajátértékét μ_{\min} . Ha

$$f_2 = \mathbf{x}^T \mathbf{Q} \mathbf{x},$$

akkor

$$f_2 \geq \mu_{\min} \mathbf{x}^T \mathbf{x}.$$

Ha

$$f = \mathbf{x}^T (\mathbf{P} + \omega \mathbf{Q}) \mathbf{x},$$

akkor

$$f \geq (\lambda_{\min} + \omega \mu_{\min}) \mathbf{x}^T \mathbf{x}.$$

Az első eset, ha

$$\mu_{\min} \geq 0.$$

Ekkor $f > 0$ minden $\omega \geq 0$ esetén. A második eset, ha

$$\mu_{\min} = -|\mu_{\min}| < 0.$$

Ekkor

$$f > (\lambda_{\min} - \omega |\mu_{\min}|) \mathbf{x}^T \mathbf{x},$$

ezért $f > 0$ minden $\omega < \frac{\lambda_{\min}}{|\mu_{\min}|}$ és $\mathbf{x} \neq \mathbf{0}$ esetén. □

Elégséges konvergencia feltételt vezethetünk le az (5.17) feltételből is [2], legyen

$$\mathbf{P} = (\mathbf{I} + \mathbf{L}^T)(\mathbf{I} + \mathbf{L}) - \mathbf{U}^T \mathbf{U},$$

$$\mathbf{Q} = (\mathbf{I} + \mathbf{U}^T)(\mathbf{I} + \mathbf{U}) - \mathbf{L}^T \mathbf{L}.$$

Ekkor $\mathbf{P} + \mathbf{Q} = \mathbf{A} + \mathbf{A}^T$. Az (5.17) feltétel alakja

$$\omega^2 \mathbf{P} - \omega(\omega - 1) \mathbf{Q} > 0. \quad (5.18)$$

Ezért a SOR módszer konvergálni fog, ha $\omega = 1$ és $\mathbf{P} > 0$. [2]

5.4. Konklúzió

Az 5.3-as tétel szerint az (5.16) és (5.17) feltételek teljesülnek, ha $\mathbf{A} + \mathbf{A}^T$ pozitív definit, és megfelelően kicsi ω -át választunk. Minden \mathbf{A} mátrix kifejezhető szimmetrikus és antiszimmetrikus mátrixok összegeként. Ezért ha az \mathbf{A} mátrixot így bontjuk fel, és a szimmetrikus mátrix pozitív definit, mindig lehetséges lesz olyan ω -t találni, hogy a SOR módszer konvergáljon. [2]

Az (5.18) feltétel más jellegű. Azt mutatja meg, hogy a SOR módszer konvergálni fog $\omega = 1$ esetén, ha

$$(\mathbf{I} + \mathbf{L}^T)(\mathbf{I} + \mathbf{L}) - \mathbf{U}^T \mathbf{U} > 0. \quad (5.19)$$

Ebből következik, hogy léteznek olyan mátrixok, amelyeknél az alsó háromszög dominánság fontos szempont a SOR módszer konvergenciájánál. A szélsőséges esetben, amikor \mathbf{U} nulla lesz, az (5.19) feltétel fennáll, és ha ω értéke 1, az egyenletek egy lépésben megoldódnak. [2]

Azt, hogy a $\mathbf{P} > 0$ és az $\mathbf{A} + \mathbf{A}^T > 0$ feltételek nem ekvivalensek, legjobban két példával lehet szemléltetni. Legyen

$$\mathbf{A} = \begin{bmatrix} 1 & -2 \\ 3 & 1 \end{bmatrix}.$$

Ebben az esetben az \mathbf{A} mátrix erősen antiszimmetrikus, és a szimmetrikus komponense pozitív definit, habár sem \mathbf{P} , sem \mathbf{Q} nem az. A SOR módszer megfelelően kicsi ω esetén konvergál. Pl. $\omega = \frac{1}{4}$. [2]

Legyen

$$\mathbf{A} = \begin{bmatrix} 1 & \frac{1}{3} \\ 2 & 1 \end{bmatrix}.$$

Az \mathbf{A} mátrix alsó háromszög domináns. A \mathbf{P} mátrix pozitív definit, de $\mathbf{A} + \mathbf{A}^T$ és \mathbf{Q} nem az. A konvergencia $\omega = 1$ esetén garantált. [2]

Érdemes kiemelni, hogy az (5.16), (5.17), (5.18) feltételek elégségesek, de nem szükségesek. A (5.16), (5.17) feltételeket nem kielégítő, nagyobb ω , nem feltétlenül jelenti, hogy a SOR módszer divergálni fog. Viszont a feltételekből látszik, hogy létezik kettő igen különböző típusa a nemszimmetrikus mátrixoknak, melyekre a SOR módszer mindig konvergál, ha megfelelő ω értéket választunk. [2]

6. A konjugált gradiens módszerek új rendszertana

6.1. A témakör felvezetése

Az egyik legjobban használható módszer az

$$\mathbf{A}\mathbf{x} = \mathbf{b} \tag{6.1}$$

egyenletrendszer megoldására, ahol \mathbf{A} valós nemszinguláris ritka mátrix, \mathbf{b} pedig valós vektor, a konjugált gradiens módszer és az ebből származtatott különböző módszerek [5]. Hestenes és Siefel eredeti 1952-es módszere [9] csak akkor alkalmazható, ha az \mathbf{A} együtthatómátrix szimmetrikus és pozitív definit. Az eredeti módszer óta már rengetek származtatott módszer született, amelyek nem csak szimmetrikus indefinit együtthatómátrixok esetén, de nemszimmetrikus együtthatómátrixok esetén is alkalmazhatóak [5]. A módszerek áttekintése azonban nehézkes lehet, például amiatt, hogy az algoritmusokat a különböző szerzők különböző módokon származtatják. Broyden 1996-os cikkében [5] rendszerezi a konjugált gradiens módszereket. Cikke támaszkodik más szerzők azonos célú korábbi munkáira, például S.F. Ashby, T.A. Manteuffel és P.E. Saylor 1990-es cikkére [1]. Ők minden algoritmust három mátrixszal jellemeznek. Az \mathbf{A} együtthatómátrixszal, egy Hermit-féle pozitív definit \mathbf{B} mátrixszal (a belső szorzat mátrixszal), és egy további \mathbf{C} mátrixszal (a prekondicionáló mátrixszal). Cikkük erősen támaszkodik Faber és Manteuffel cikkére [8], akik a valós \mathbf{M} mátrixot B-normális(s)-nek nevezik, ha a \mathbf{B} mátrix Hermit-féle, pozitív definit, és

$$\mathbf{M}^T \mathbf{B} = \mathbf{B} p(\mathbf{M}), \tag{6.2}$$

ahol s a legkisebb fokszáma a $p(\mathbf{M})$ mátrixpolinomnak, amelyre a (6.2) egyenlőség teljesül [8]. Megmutatták, hogy egy 3-tagú rekurrens prekondicionált konjugált gradiens módszernél, ha a \mathbf{CA} mátrixszorzat B-normális(1), akkor a módszer véges lépésben leáll [8]. Broyden rendszertanában egy \mathbf{G} Hesse-mátrixból és egy további \mathbf{K} mátrixból származtatja a módszereket [5]. Broyden \mathbf{G} mátrixa általában azonos az Ashby-féle \mathbf{B} mátrixszal, azzal a különbséggel, hogy a pozitív definitég a \mathbf{G}

mátrixnál nem megkövetelt [5]. További összefüggés a két rendszertan között, hogy

$$\mathbf{C}\mathbf{A} = \mathbf{K}\mathbf{G} \quad (6.3)$$

ahol \mathbf{C} és \mathbf{A} az Ashby féle \mathbf{C} és \mathbf{A} mátrix [5]. Ebből következik, hogy a $\mathbf{K}\mathbf{G}$ mátrixszorzatnak B-normális(1)-nek kell lennie. Legyen $\mathbf{K}\mathbf{G}$ az \mathbf{M} mátrix, és \mathbf{G} a \mathbf{B} mátrix a (6.2) egyenlőségben. Ekkor a B-normalitás követelménye

$$\mathbf{G}\mathbf{K}^T\mathbf{G} = \mathbf{G}p(\mathbf{K}\mathbf{G}), \quad (6.4)$$

ahol $p(\mathbf{K}\mathbf{G})$ lineáris mátrixpolinom [5]. Ennek elégséges feltétele, ha a \mathbf{K} mátrix szimmetrikus, ekkor $p(\mathbf{K}\mathbf{G})$ helyett vehetjük a $\mathbf{K}\mathbf{G}$ mátrixszorzatot [5]. Habár a \mathbf{K} mátrix szimmetrikussága erősebb feltétel, mint a B-normalitás, de ahhoz még elég általános, hogy a standard módszereket is tartalmazza a rendszertan [5]. Ennek a mátrixfelbontásnak még egy következménye, hogy a hagyományos 2-tagú konjugált gradiens módszereknél, a két tipikus numerikus instabilitása ezeknek a módszereknek (lásd később), a \mathbf{G} és a \mathbf{K} mátrixokból könnyen megállapítható [5]. A \mathbf{K} és \mathbf{G} mátrixokon kívül még szükségünk van bizonyos vektorokra, a generátor vektorokra, amelyek minden iterációval változnak. Ezeket kétféleképpen választhatjuk, az egyik választás vezet a 2-tagú (Hestenes-Stiefel) rekurrens módszerekhez, míg a másik választás a 3-tagú rekurrens (Lánczos) módszerekhez [5].

Legyen a $\phi(\mathbf{x})$ kvadratikus függvény

$$\phi(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T\mathbf{G}\mathbf{x} - \mathbf{x}^T\mathbf{h}, \quad (6.5)$$

ahol \mathbf{G} n -ed rendű szimmetrikus mátrix, de nem feltétlen pozitív definit, \mathbf{h} pedig n -ed rendű vektor. Legyen \mathbf{x}_1 a kezdeti értéke az \mathbf{x} vektornak, és legyen \mathbf{S}_i $n \times i$ -s, vagy $n \times 2i$ -s mátrix. Ekkor, ha \mathbf{x}_{i+1} jelöli azt az \mathbf{x} vektor értéket, amelyre a $\phi(\mathbf{x})$ függvény stacionárius az

$$\mathbf{x} = \mathbf{x}_1 + \mathbf{S}_i\mathbf{z} \quad (6.6)$$

hipersík felett, ahol \mathbf{z} i -ed vagy $2i$ -ed rendű vektor, akkor

$$\mathbf{x}_{i+1} = \mathbf{x}_1 - \mathbf{S}_i(\mathbf{S}_i^T\mathbf{G}\mathbf{S}_i)^{-1}\mathbf{S}_i^T\mathbf{g}_1, \quad (6.7)$$

ahol \mathbf{g}_1 a gradiensvektora a ϕ függvénynek az \mathbf{x}_1 helyen. Továbbá, ha

$$\mathbf{S}_i = \begin{bmatrix} \mathbf{P}_1 & \mathbf{P}_2 & \dots & \mathbf{P}_i \end{bmatrix}, \quad (6.8)$$

ahol \mathbf{P}_j , $1 \leq j \leq i$, részmátrixok, egy vagy két oszloppal, úgy, hogy a $\mathbf{C}_j = \mathbf{P}_j^T\mathbf{G}\mathbf{P}_j$

mátrix nemszinguláris, és

$$\mathbf{P}_j^T \mathbf{G} \mathbf{P}_k = \mathbf{0}, \quad (6.9)$$

ahol $j \neq k$, $1 \leq j, k \leq i$. Az \mathbf{x}_{i+1} vektor kifejezhető úgy is, hogy

$$\mathbf{x}_{i+1} = \mathbf{x}_i - \mathbf{P}_i (\mathbf{P}_i^T \mathbf{G} \mathbf{P}_i)^{-1} \mathbf{P}_i^T \mathbf{g}_i, \quad (6.10)$$

ahol

$$\mathbf{g}_i = \mathbf{G} \mathbf{x}_i - \mathbf{h}. \quad (6.11)$$

A (6.7) egyenlőséget a \mathbf{G} mátrixszal szorozva, és a \mathbf{h} vektort mindkét oldalból kivonva,

$$\mathbf{g}_{i+1} = \mathbf{Q}_i \mathbf{g}_1, \quad (6.12)$$

ahol

$$\mathbf{Q}_i = \mathbf{I} - \mathbf{G} \mathbf{S}_i (\mathbf{S}_i^T \mathbf{G} \mathbf{S}_i)^{-1} \mathbf{S}_i^T \quad (6.13)$$

[5].

Ahhoz, hogy ki tudjuk számolni az \mathbf{x}_{i+2} vektort, szükség van a \mathbf{P}_{i+1} mátrixra. Könnyen látható, hogy ha

$$\mathbf{P}_{i+1} = \mathbf{Q}_i^T \mathbf{W}_{i+1}, \quad (6.14)$$

akkor a (6.9) egyenlőség fennáll $(i+1)$ -re, tetszőleges \mathbf{W}_{i+1} mátrixszal. A \mathbf{Q}_i mátrix túl nagy és sűrű, de mivel a \mathbf{W}_{i+1} mátrix tetszőlegesen választható, ezzel egyszerűsíthetjük a (6.14) összefüggést. Két lehetséges választása a \mathbf{W}_{i+1} mátrixnak különösen előnyös. [5]

Az első lehetőség, amikor a \mathbf{G} mátrixnak nincs semmilyen különösebb blokk struktúrája, és amikor a \mathbf{P}_j mátrixok \mathbf{p}_j vektorok, a kapcsolódó \mathbf{W}_j generátor mátrixok pedig \mathbf{w}_j vektorok. Megmutatjuk, hogy ha

$$\mathbf{w}_j = \mathbf{K} \mathbf{g}_j, \quad (6.15)$$

ahol \mathbf{K} tetszőleges szimmetrikus mátrix, akkor a (6.14) egyenlőség a

$$\mathbf{p}_{i+1} = \mathbf{K} \mathbf{g}_{i+1} - \mathbf{p}_i \alpha_i \quad (6.16)$$

alakra egyszerűsödik, ahol α_i konstans. [5]

Először bebizonyítjuk a következő tételt.

6.1. Tétel. *Legyen \mathbf{x}_{i+1} stacionárius pontja a $\phi(\mathbf{x})$ függvénynek a (6.6) egyenlőségben definiált hipersík felett. Legyen \mathbf{g}_{i+1} az \mathbf{x}_{i+1} ponthoz tartozó gradiense a $\phi(\mathbf{x})$ függvénynek, és legyenek a \mathbf{w}_j vektorok az \mathbf{S}_i mátrix generátorai. Ekkor $\mathbf{w}_j^T \mathbf{g}_{i+1} = 0$, $i \leq j \leq i$. [5, 9. oldal, Lemma 1]*

Bizonyítás. A (6.12) egyenlőségből $\mathbf{w}_j^T \mathbf{g}_{i+1} = \mathbf{w}_j^T \mathbf{Q}_i \mathbf{g}_1$, a (6.13) és a (6.9) egyenlő-

ségből $\mathbf{Q}_{j-1}\mathbf{Q}_i = \mathbf{Q}_i$, $0 \leq j-1 \leq i$, ezért $\mathbf{w}_j^T \mathbf{g}_{i+1} = \mathbf{w}_j^T \mathbf{Q}_{j-1} \mathbf{Q}_i \mathbf{g}_1$. Ezért a (6.14) egyenlőségből $\mathbf{w}_j^T \mathbf{g}_{i+1} = \mathbf{p}_j^T \mathbf{Q}_i \mathbf{g}_1$. A tétel következik a (6.8) egyenlőségből, mivel a (6.13) egyenlőség miatt $\mathbf{S}_i^T \mathbf{Q}_i = 0$. \square

Következmény. Ha $\mathbf{w}_j = \mathbf{K} \mathbf{g}_j$, $1 \leq j \leq i$, akkor $\mathbf{g}_j^T \mathbf{K} \mathbf{g}_k = 0$, $j \neq k$, $j \geq 1$, $k \leq i+1$.

Megmutatjuk, hogy ha a \mathbf{K} mátrix definit, akkor a generátor vektorok fenti választásával

$$\mathbf{p}_j^T \mathbf{G} \mathbf{K} \mathbf{g}_{i+1} = 0, \quad (6.17)$$

ahol $1 \leq j \leq i-1$. A (6.10) egyenlőségben i -t j -re cserélve, az egyenlőséget megszorozva a $\mathbf{K} \mathbf{G}$ mátrixszal, és a $\mathbf{K} \mathbf{h}$ vektort kivonva mindkét oldalból, a (6.11) egyenlőségből kapjuk, hogy

$$\mathbf{K} \mathbf{g}_{j+1} = \mathbf{K} \mathbf{g}_j - \mathbf{K} \mathbf{G} \mathbf{p}_j \gamma_j, \quad (6.18)$$

ahol $\gamma_j = \mathbf{p}_j^T \mathbf{g}_j / \mathbf{p}_j^T \mathbf{G} \mathbf{p}_j$. A \mathbf{P}_j mátrix most egy vektor, ezért \mathbf{p}_j -vel jelöltük. Megszorozva ezt az egyenlőséget a \mathbf{g}_{i+1}^T vektorral, a 6.1-es tétel következményéből kapjuk, hogy $\gamma_j \mathbf{p}_j^T \mathbf{G} \mathbf{K} \mathbf{g}_{i+1} = 0$, $1 \leq j \leq i-1$. A (6.18) egyenlőséget megszorozva a \mathbf{g}_{j+1}^T vektorral, a 6.1-es tétel következménye, hogy ha $\mathbf{g}_{j+1} \neq \mathbf{0}$, és minthogy \mathbf{K} definit, akkor $\gamma \neq 0$, így kapjuk a (6.17) egyenlőséget [5].

A (6.13) és a (6.14) egyenlőségből

$$\mathbf{p}_{i+1} = \mathbf{K} \mathbf{g}_{i+1} - \mathbf{S}_i (\mathbf{S}_i^T \mathbf{G} \mathbf{S}_i)^{-1} \mathbf{S}_i^T \mathbf{G} \mathbf{K} \mathbf{g}_{i+1}, \quad (6.19)$$

de a (6.17) egyenlőség miatt ennek az egyenlőségnek csak a $\mathbf{S}_i^T \mathbf{G} \mathbf{K} \mathbf{g}_{i+1}$ tagja nem nulla, ebből következik a (6.16) egyenlőség és az $\mathbf{S}_i^T \mathbf{G} \mathbf{S}_i$ mátrix diagonalitása [5].

Az α_i konstansot úgy választjuk a (6.16) egyenlőségben, hogy $\mathbf{p}_i^T \mathbf{G} \mathbf{p}_{i+1} = 0$ legyen. Ez mindig lehetséges, ha $\mathbf{p}_i^T \mathbf{G} \mathbf{p}_i \neq 0$. Ez a feltétel mindig teljesül a nem nulla \mathbf{p}_i vektorokra, ha a \mathbf{G} mátrix pozitív definit. Az olyan algoritmusokat, ahol a \mathbf{G} mátrix pozitív definit, b-stabilisnak nevezzük. Ha a \mathbf{K} mátrix definit, egy másféle numerikus instabilitást kerülhetünk el, amikor is az egymást követő lépések túl kicsik lesznek [5]. Azokat az algoritmusokat, ahol a \mathbf{K} mátrix definit, ω -stabilisnak nevezzük. Így a konjugált gradiens algoritmusok kétféle numerikus instabilitása a \mathbf{G} és a \mathbf{K} mátrixoktól függenek [5].

A második lehetőség a \mathbf{w}_j vektor választására, ha

$$\mathbf{w}_j = \mathbf{K} \mathbf{G} \mathbf{p}_{j-1}, \quad (6.20)$$

ahol a \mathbf{p}_1 vektort tetszőlegesen választhatjuk és $\mathbf{p}_0 = \mathbf{0}$. Ezzel a választással a (6.14)

egyenlőség alakja

$$\mathbf{p}_{i+1} = \mathbf{K}\mathbf{G}\mathbf{p}_i - \mathbf{p}_i\alpha_i - \mathbf{p}_{i-1}\beta_{i-1}, \quad (6.21)$$

ahol az α_i és a β_{i-1} konstansokat úgy választjuk, hogy

$$\mathbf{p}_{i-1}^T \mathbf{G}\mathbf{p}_{i+1} = \mathbf{p}_i^T \mathbf{G}\mathbf{p}_{i+1} = 0. \quad (6.22)$$

6.2. Tétel. *Legyenek a \mathbf{p}_j vektorok a (6.14) és (6.13) egyenlőségben definiált \mathbf{w}_j generátorvektorokból számolva, $1 \leq j \leq i$. Ekkor*

$$\mathbf{w}_j^T \mathbf{G}\mathbf{p}_i = 0, \quad (6.23)$$

ahol $1 \leq j \leq i-1$. [5, 10. oldal, Lemma 2]

Bizonyítás. A (6.8) egyenlőség és a \mathbf{p}_j vektorok konjugáltsága miatt, $\mathbf{S}_j^T \mathbf{G}\mathbf{p}_i = 0$, $1 \leq j \leq i-1$. A (6.13)-as egyenlőségből, mivel $\mathbf{Q}_0 = \mathbf{I}$, $\mathbf{G}\mathbf{p}_i = \mathbf{Q}_{j-1} \mathbf{G}\mathbf{p}_i$, $1 \leq j \leq i-1$. Ezért $\mathbf{w}_j^T \mathbf{G}\mathbf{p}_i = \mathbf{w}_j^T \mathbf{Q}_{j-1} \mathbf{G}\mathbf{p}_i$, és a (6.14) egyenlőségből és a konjugáltságból, $\mathbf{w}_j^T \mathbf{G}\mathbf{p}_i = \mathbf{p}_j^T \mathbf{G}\mathbf{p}_i = 0$, $1 \leq j \leq i-1$. \square

Következmény. *Ha a \mathbf{w}_j generátorvektorokat a (6.20) egyenlőség szerint választjuk, akkor*

$$\mathbf{p}_j^T \mathbf{G}\mathbf{K}\mathbf{G}\mathbf{p}_i = 0,$$

ahol $0 \leq j \leq i-2$.

A generátor vektorok fenti választásával, a (6.14) egyenlőség alakja

$$\mathbf{p}_{i+1} = \mathbf{K}\mathbf{G}\mathbf{p}_i - \mathbf{S}_i(\mathbf{S}_i^T \mathbf{G}\mathbf{S}_i)^{-1} \mathbf{S}_i^T \mathbf{G}\mathbf{K}\mathbf{G}\mathbf{p}_i. \quad (6.24)$$

A (6.24) egyenlőségből, csak az utolsó két tagja nem nulla a $\mathbf{S}_i^T \mathbf{G}\mathbf{K}\mathbf{G}\mathbf{p}_i$ szorzatnak, ezért a (6.21) egyenlőség az $\mathbf{S}_i^T \mathbf{G}\mathbf{S}_i$ mátrix diagonalitásából azonnal következik [5].

Néhány megjegyzés:

1. A kezdeti \mathbf{p}_1 vektor tetszőleges. Ha a \mathbf{p}_1 vektornak a \mathbf{g}_1 vektort választjuk, akkor az orthodir algoritmust kapjuk [1] [5].
2. Ha $\mathbf{p}_1 = \mathbf{K}\mathbf{g}_1$, akkor a (6.20) egyenlőségből kapott $\{\mathbf{x}_i\}$ vektorsorozat megegyezik a (6.16) egyenlőségből kapott $\{\mathbf{x}_i\}$ vektorsorozattal, és a stabilitás feltételei hasonlóak [5].
3. Mivel a kezdeti \mathbf{p}_1 vektor választása tetszőleges, választhatjuk a \mathbf{P}_1 mátrixnak, tetszőleges számú oszloppal, és így kapjuk a blokk módszereket [5].
4. A 2-tagú módszereknél a konjugált \mathbf{p}_j vektorok számítása elválaszthatatlanul összefügg a $\phi(\mathbf{x})$ függvény stacionárius pontjának számításával, viszont a 3-

tagú módszereknél a \mathbf{p}_j konjugált vektorok generálása teljesen független a $\phi(\mathbf{x})$ függvény stacionárius pontjának számításától (módszer 6, 16, 17, 19) [5].

6.2. Szimmetrikus mátrixok

A fő módszerek az (6.1) egyenletrendszer megoldására, ha az \mathbf{A} együtthatómátrix szimmetrikus, a következők [5].

Módszer	\mathbf{G}	\mathbf{h}	\mathbf{K}	Név	Referencia
1	\mathbf{A}	\mathbf{b}	\mathbf{I}	cg	[7, 3.5.1-es fejezet], [9]
2	\mathbf{A}^2	\mathbf{Ab}	\mathbf{A}^{-1}	cr	[7, 3.6.1-es fejezet], [9]
3	\mathbf{A}	\mathbf{b}	\mathbf{M}^{-1}	pcg	[1]

Ezek a \mathbf{p}_i vektorra 2-tagú rekurrens módszert adnak. Az első kettő, az eredeti konjugált gradiens módszer és a konjugált reziduális módszer b-stabilis és ω -stabilis, ha az \mathbf{A} együtthatómátrix pozitív definit. Ha az \mathbf{A} együtthatómátrix indefinit, akkor a cg módszer csak \mathbf{b} -stabilis, a cr módszer pedig csak ω -stabilis. Ezért a módszerek használata indefinit \mathbf{A} együtthatómátrix esetén nem ajánlott. A 3-as módszer, a prekondicionált konjugált gradiens módszer b-stabilis és ω -stabilis, ha az \mathbf{A} és \mathbf{M} mátrixok definitek. Az \mathbf{M} prekondicionáló mátrixot úgy választjuk, hogy a \mathbf{KG} mátrix kondíciós számát csökkentse, ezzel javítva a módszer konvergenciáját. A fejezet több másik algoritmusát is prekondicionálhatjuk a \mathbf{K} mátrix módosításával, de erre külön nem térünk ki [5].

A \mathbf{p}_i vektorra 3-tagú rekurrens módszert adnak a következők, ha az \mathbf{A} együtthatómátrix szimmetrikus [5].

Módszer	\mathbf{G}	\mathbf{h}	\mathbf{K}	Név	Referencia
4	\mathbf{A}	\mathbf{b}	\mathbf{I}	Nazareth	[10]
5	\mathbf{A}^2	\mathbf{Ab}	\mathbf{A}^{-1}		[5]
6	\mathbf{I}	$\mathbf{A}^{-1}\mathbf{b}$	\mathbf{A}	SYMMLQ	[7, 3.8.1-es fejezet]

A 4-es és 5-ös módszer az 1-es és 2-es módszer 3-tagú változata. A 6-os módszernek nincs 2-tagú megfelelője. Mivel $\mathbf{G} = \mathbf{I}$, bármilyen ortogonális \mathbf{p}_j vektorok megfelelőek lesznek. A probléma az \mathbf{x}_{i+1} vektorok kiszámítása, mert a (6.7) egyenlőségben szerepel $\mathbf{g}_i = \mathbf{x}_1 - \mathbf{A}^{-1}\mathbf{b}$, és \mathbf{A}^{-1} nem ismert. Egy lehetséges megoldás, hogy válasszuk a \mathbf{p}_1 vektort úgy, hogy $\mathbf{p}_1 = \mathbf{A}\mathbf{r}_1$, ahol

$$\mathbf{r}_i = \mathbf{A}\mathbf{x}_i - \mathbf{b} \quad (6.25)$$

maradéktag, és indirekt kapjuk a (6.10) egyenlőségben szereplő $\mathbf{p}_i^T \mathbf{g}_i$ -t, de ez numerikusan nemkívánatos [5]. Ha a \mathbf{p}_j vektorokat a (6.21) egyenlőség felhasználásával

számoljuk, akkor a (6.8) egyenlőségből, és a \mathbf{G} és \mathbf{K} mátrixok megfelelő értékeit felhasználva,

$$\mathbf{A}\mathbf{S}_i = \mathbf{S}_{i+1}\mathbf{T}_{i+1}, \quad (6.26)$$

ahol \mathbf{T}_{i+1} az $(i+1) \times i$ méretű bal felső sarokminormátrixa valamilyen tridiagonális mátrixnak [5]. Legyen \mathbf{Q}_{i+1} olyan ortogonális mátrix, hogy

$$\mathbf{Q}_{i+1}\mathbf{T}_{i+1} = \begin{bmatrix} \mathbf{U}_i \\ \mathbf{0}^T \end{bmatrix}, \quad (6.27)$$

ahol \mathbf{U}_i felső háromszögmátrix, és $\mathbf{0}^T$ sorvektor. Legyen az \mathbf{Y}_i mátrix

$$\begin{bmatrix} \mathbf{Y}_i & \mathbf{y}_{i+1} \end{bmatrix} = \mathbf{S}_{i+1}\mathbf{Q}_{i+1}^T. \quad (6.28)$$

Ekkor a (6.26), (6.27), (6.28) egyenlőségekből

$$\mathbf{A}\mathbf{S}_i = \mathbf{Y}_i\mathbf{U}_i, \quad (6.29)$$

ahol, mivel az \mathbf{S}_{i+1} mátrix oszlopai ortogonálisak, az \mathbf{Y}_i mátrix oszlopai szintén ortogonálisak. Ha most minimalizálni akarjuk a $\phi(\mathbf{x})$ függvényt az

$$\mathbf{x} = \mathbf{x}_1 + \mathbf{Y}_i\mathbf{z} \quad (6.30)$$

hipersík felett, a \mathbf{G}_i , \mathbf{S}_i mátrix és a \mathbf{g}_i vektor megfelelő értékeit a (6.7) egyenlőségbe helyettesítve, a (6.11), (6.25) és (6.29) egyenlőségből kapjuk, hogy

$$\mathbf{x}_{i+1} = \mathbf{x}_i - \mathbf{Y}_i(\mathbf{Y}_i^T\mathbf{Y}_i)^{-1}\mathbf{U}_i^{-T}\mathbf{S}_i^T\mathbf{r}_1. \quad (6.31)$$

Ez lényegében az SYMMLQ módszere Paige-nek és Saunders-nek [5]. Ők azonban az \mathbf{x}_i vektorokat csak segédvektoroknak használták, és az \mathbf{x}_i vektorokból számoltak egy közelítő vektorsorozatot az (6.1) egyenletrendszer megoldására, számos más fontos numerikus finomítás mellett. Az algoritmusukat először Fletcher használta mint minimum-hiba algoritmus [5], illetve mivel szükség van a \mathbf{p}_{i+1} vektorra az \mathbf{x}_{i+1} vektor kiszámításához, tekinthetünk a módszerre mint egy implicit „előrettekintő” módszerre [5].

6.3. Nemszimmetrikus mátrixok

Ha az \mathbf{A} együtthatómátrix nemszimmetrikus, az 1-6 módszereket nem tudjuk alkalmazni. Az 1-es módszert adja azonban a 7-es módszert és még két másik variációt. [5].

Módszer	\mathbf{G}	\mathbf{h}	\mathbf{K}	Név	Referencia
7	$\mathbf{A}^T \mathbf{A}$	$\mathbf{A}^T \mathbf{b}$	\mathbf{I}	cgne	[7, 3.7.3-as fejezet], [9]
8	\mathbf{I}	$\mathbf{A}^{-1} \mathbf{b}$	$\mathbf{A}^T \mathbf{A}$	Craig's method	[7, 3.7.3-as fejezet]
9	$\mathbf{Z} \mathbf{A}$	$\mathbf{Z} \mathbf{b}$	\mathbf{Z}^{-1}	orthodir	[7, 2.5-ös fejezet]

A 7-es módszer, annak ellenére, hogy b-stabilis és ω -stabilis, általában nem kielégítő, mert a \mathbf{GK} mátrix kondíciószáma a négyzete az 1-6 módszerek \mathbf{GK} mátrixának kondíciószámának, és emiatt gyakran lassú a konvergencia. A 8-as módszerre hasonlóak igazak, de ennél a módszernél a 6-os módszerhez hasonlóan az \mathbf{A}^{-1} mátrixra is szükség van még, az \mathbf{x}_i vektorok kiszámításához. Ezen túlléphetünk, ha megszorozzuk balról a (6.16) egyenlőséget az \mathbf{A}^{-T} mátrixszal, így nem a $\{\mathbf{p}_i\}$ vektorsorozatot, hanem egy $\{\mathbf{q}_i\}$ vektorsorozatot generálva, ahol $\mathbf{q}_i = \mathbf{A}^{-T} \mathbf{p}_i$, amikből az $\{\mathbf{x}_i\}$ sorozat számolható [5]. A 9-es módszer egy speciális esete az általánosabb Young és Jea módszernek [5], ahol szükséges, hogy \mathbf{Z} és $\mathbf{Z} \mathbf{A}$ mátrixok szimmetrikusak legyenek, ami egy erős megkötés. A stabilitás garantált, ha \mathbf{Z} és $\mathbf{Z} \mathbf{A}$ is definit [5].

Mivel a fenti módszerek egyike se teljesen kielégítő nemszimmetrikus mátrixok esetén, nézzük meg a speciális esetet, amikor a \mathbf{G} mátrix alakja

$$\mathbf{G} = \begin{bmatrix} \mathbf{G}_{11} & \mathbf{0} \\ \mathbf{0} & \mathbf{G}_{22} \end{bmatrix}, \quad (6.32)$$

vagy

$$\mathbf{G} = \begin{bmatrix} \mathbf{0} & \mathbf{G}_{12} \\ \mathbf{G}_{21} & \mathbf{0} \end{bmatrix}. \quad (6.33)$$

Innen nem nehéz megmutatni [5], hogy ha a $2n \times 2$ -es \mathbf{W}_j mátrix alakja

$$\mathbf{W}_j = \begin{bmatrix} \mathbf{w}_{j1} & \mathbf{0} \\ \mathbf{0} & \mathbf{w}_{j2} \end{bmatrix}, \quad j = 1, 2, \dots, i, \quad (6.34)$$

akkor

$$\mathbf{Q}_i = \begin{bmatrix} \mathbf{Q}_{i1} & \mathbf{0} \\ \mathbf{0} & \mathbf{Q}_{i2} \end{bmatrix} \quad (6.35)$$

alakú, és

$$\mathbf{P}_{i+1} = \begin{bmatrix} \mathbf{u}_{i+1} & \mathbf{0} \\ \mathbf{0} & \mathbf{v}_{i+1} \end{bmatrix} \quad (6.36)$$

alakú. Így a (6.14) egyenlőségből

$$\mathbf{u}_{i+1} = \mathbf{Q}_{i1}^T \mathbf{w}_{i+1,1} \quad (6.37)$$

[5]. Hasonlóan kaphatjuk a \mathbf{v}_{i+1} vektort.

Úgy, mint a (6.14) egyenlőség esetében, a \mathbf{W}_{i+1} mátrixot megfelelően kell választani, hogy működő algoritmust kapjunk, és itt is két eset lehetséges [5]. Az első, amikor

$$\mathbf{W}_j = \mathbf{K} \mathbf{F}_j \quad (6.38)$$

alakú, ahol

$$\mathbf{F}_j = \begin{bmatrix} \mathbf{g}_{j1} & \mathbf{0} \\ \mathbf{0} & \mathbf{g}_{j2} \end{bmatrix}. \quad (6.39)$$

A \mathbf{g}_{j1} és \mathbf{g}_{j2} vektorok egyértelmű partíciói a gradiens \mathbf{g}_j vektornak a (6.11), (6.32) és (6.33) egyenlőségek miatt. A \mathbf{K} mátrix alakja

$$\mathbf{K} = \begin{bmatrix} \mathbf{K}_{11} & \mathbf{0} \\ \mathbf{0} & \mathbf{K}_{22} \end{bmatrix}, \quad (6.40)$$

vagy

$$\mathbf{K} = \begin{bmatrix} \mathbf{0} & \mathbf{K}_{12} \\ \mathbf{K}_{21} & \mathbf{0} \end{bmatrix}, \quad (6.41)$$

ahol \mathbf{K}_{11} , \mathbf{K}_{12} , \mathbf{K}_{21} , \mathbf{K}_{22} konstansok és $\mathbf{K} = \mathbf{K}^T$. Ebben az esetben meg lehet mutatni [5], hogy a (6.37) egyenlőségből

$$\mathbf{u}_{i+1} = \mathbf{w}_{i+1,1} - \mathbf{u}_i \alpha_{i1}, \quad (6.42)$$

ahol α_{i1} -et úgy választjuk, hogy kielégítse vagy az $\mathbf{u}_i^T \mathbf{G}_{11} \mathbf{u}_{i+1} = 0$ egyenlőséget, ha a \mathbf{G} mátrix (6.32) szerint adott, vagy a $\mathbf{v}_i^T \mathbf{G}_{21} \mathbf{u}_{i+1} = 0$ egyenlőséget, ha a \mathbf{G} mátrix (6.33) szerint adott [5]. Hasonlóan kaphatjuk a \mathbf{v}_{i+1} vektort. A \mathbf{G} és a \mathbf{K} mátrix lehetséges választásai jelenleg a következők [5].

Módszer	\mathbf{G}	\mathbf{h}	\mathbf{K}	Név	Referencia
10	$\begin{bmatrix} \mathbf{0} & \mathbf{A}^T \\ \mathbf{A} & \mathbf{0} \end{bmatrix}$	$\begin{bmatrix} \mathbf{c} \\ \mathbf{b} \end{bmatrix}$	$\begin{bmatrix} \mathbf{0} & \mathbf{I} \\ \mathbf{I} & \mathbf{0} \end{bmatrix}$	bcg	[7, 3.5.2-es fejezet]
11	$\begin{bmatrix} \mathbf{0} & \mathbf{A}^T \\ \mathbf{A} & \mathbf{0} \end{bmatrix}$	$\begin{bmatrix} \mathbf{c} \\ \mathbf{b} \end{bmatrix}$	$\begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}$	Heg	[7, 3.5.4-es fejezet]
12	$\begin{bmatrix} \mathbf{A}^T \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{A} \mathbf{A}^T \end{bmatrix}$	$\begin{bmatrix} \mathbf{A}^T \mathbf{b} \\ \mathbf{A} \mathbf{c} \end{bmatrix}$	$\begin{bmatrix} \mathbf{0} & \mathbf{A}^{-1} \\ \mathbf{A}^{-T} & \mathbf{0} \end{bmatrix}$	bcr	[5]

Az \mathbf{A}^{-T} jelölés jelentése $(\mathbf{A}^{-1})^T$. Ezeknél a módszereknél, mint az 1-es és 2-es módszereknél, ha a \mathbf{G} és a \mathbf{K} mátrix nem definit, a stabilitás nem garantált. A 10-es módszer a Lanczos-féle bikonjugált gradiens módszer, ami se nem b-stabilis, se nem ω -stabilis. A 11-es Hegedüs-féle módszer ω -stabilis, de nem b-stabilis. A 12-es mód-

szer a 2-es módszer általánosítása, és b-stabilis, ha az \mathbf{A} együtthatómátrix négyzetes és nonszinguláris [5].

Az algoritmusok következő csoportját úgy kapjuk, hogy a (6.38) egyenlőséget a

$$\mathbf{W}_j = \mathbf{KGP}_{j-1} \quad (6.43)$$

egyenlőségre cseréljük [5]. Ekkor meg lehet mutatni [5], hogy

$$\mathbf{u}_{i+1} = \mathbf{w}_{i+1,1} - \mathbf{u}_i \alpha_{i1} - \mathbf{u}_{i-1} \beta_{i-1,1}, \quad (6.44)$$

ahol α_{i1} -et és $\beta_{i-1,1}$ -et úgy választjuk, hogy kielégítse az

$$\mathbf{u}_{i-1}^T \mathbf{G}_{11} \mathbf{u}_{i+1} = \mathbf{u}_i^T \mathbf{G}_{11} \mathbf{u}_{i+1} = 0$$

egyenlőséget, ha a \mathbf{G} mátrix (6.32) szerint adott, vagy a

$$\mathbf{v}_{i-1}^T \mathbf{G}_{21} \mathbf{u}_{i+1} = \mathbf{v}_i^T \mathbf{G}_{21} \mathbf{u}_{i+1} = 0$$

egyenlőséget, ha a \mathbf{G} mátrix (6.33) szerint adott. Úgy, mint a 4-6-os módszereknél, \mathbf{u}_0 nullvektor és \mathbf{u}_1 tetszőleges vektor. A (6.44) egyenlőséghez hasonló egyenlőségből kaphatjuk a \mathbf{v}_{i+1} vektort. A következők a lehetőségek [5].

Módszer	\mathbf{G}	\mathbf{h}	\mathbf{K}	Név	Referencia
13	$\begin{bmatrix} \mathbf{0} & \mathbf{A}^T \\ \mathbf{A} & \mathbf{0} \end{bmatrix}$	$\begin{bmatrix} \mathbf{c} \\ \mathbf{b} \end{bmatrix}$	$\begin{bmatrix} \mathbf{0} & \mathbf{I} \\ \mathbf{I} & \mathbf{0} \end{bmatrix}$	MRZ	[7, 3.5.3-as fejezet]
14	$\begin{bmatrix} \mathbf{0} & \mathbf{A}^T \\ \mathbf{A} & \mathbf{0} \end{bmatrix}$	$\begin{bmatrix} \mathbf{c} \\ \mathbf{b} \end{bmatrix}$	$\begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}$	Heg3	[7, 3.5.3-as fejezet]
15	$\begin{bmatrix} \mathbf{A}^T \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{A} \mathbf{A}^T \end{bmatrix}$	$\begin{bmatrix} \mathbf{A}^T \mathbf{b} \\ \mathbf{A} \mathbf{c} \end{bmatrix}$	$\begin{bmatrix} \mathbf{0} & \mathbf{A}^{-1} \\ \mathbf{A}^{-T} & \mathbf{0} \end{bmatrix}$	bcr3	[7, 3.6.4-es fejezet]
16	$\begin{bmatrix} \mathbf{0} & \mathbf{I} \\ \mathbf{I} & \mathbf{0} \end{bmatrix}$	$\begin{bmatrix} \mathbf{A}^{-T} \mathbf{c} \\ \mathbf{A}^{-1} \mathbf{b} \end{bmatrix}$	$\begin{bmatrix} \mathbf{0} & \mathbf{A} \\ \mathbf{A}^T & \mathbf{0} \end{bmatrix}$	QMR	[7, 1.2-es fejezet]
17	$\begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}$	$\begin{bmatrix} \mathbf{A}^{-1} \mathbf{b} \\ \mathbf{A}^{-T} \mathbf{c} \end{bmatrix}$	$\begin{bmatrix} \mathbf{0} & \mathbf{A}^T \\ \mathbf{A} & \mathbf{0} \end{bmatrix}$		[5]

A 13-as módszer a 3-tagú változata a bikonjugált gradiens módszernek. A 14-es és 15-ös módszer a 3-tagú változata a Hegedüs és bikonjugált reziduális módszernek. A 15-ös módszer érdekes tulajdonsága, hogy az \mathbf{A} mátrix nem kell, hogy négyzetes legyen, és ha $m \times n$ -es n rangú mátrix (tehát $\mathbf{A}^T \mathbf{A}$ nonszinguláris), akkor a \mathbf{K}_{12} mátrixot $(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$ mátrixnak választva, az algoritmus változatlan és használható. A 14-es módszernek ugyanez a tulajdonsága. A 16-os módszer alkotja az alapját

a QMR módszernek. Csak úgy, mint a SYMMLQ (módszer 6) esetén, a probléma a kvadrátikus függvény stacionárius pontjának számolása, mert szükség van \mathbf{A}^{-1} -re. A 17-es módszereknél ugyanolyan nehézségekbe ütközünk, mint a 6-os és 16-os módszernél, és ezeket hasonlóan is lehet megoldani [5].

A módszerek utolsó csoportja tulajdonképpen az általánosított Lanczos-féle algoritmus [5]. Legyen a (6.21) egyenlőségben a \mathbf{G} és \mathbf{K} mátrix a (6.32) és a (6.41) egyenlőség szerint adott. Legyen $\mathbf{p}_i^T = [\mathbf{u}_i^T \ \mathbf{v}_i^T]$. Ha a tetszőlegesen választható $\mathbf{p}_1 = \mathbf{0}$ és $\mathbf{v}_1 = \mathbf{0}$, akkor a (6.21) egyenlőségből kivonva, és a konjugáltságot megkövetelve kapjuk, hogy

$$\mathbf{u}_{i+1} = \mathbf{K}_{12}\mathbf{G}_{22}\mathbf{v}_i - \mathbf{u}_{i-1}\beta_{i-1} \quad (6.45)$$

és $\mathbf{v}_{i+1} = \mathbf{0}$, minden páros i -re, és

$$\mathbf{v}_{i+1} = \mathbf{K}_{21}\mathbf{G}_{11}\mathbf{u}_i - \mathbf{v}_{i-1}\delta_{i-1} \quad (6.46)$$

és $\mathbf{u}_{i+1} = \mathbf{0}$, minden páratlan i -re [5]. Hasonló egyenlőségekhez juthatunk a $\mathbf{p}_1 = \mathbf{0}$ és az $\mathbf{u}_1 = \mathbf{0}$ választással [5]. Ez a Golub-Kahan algoritmus általánosítása, és két esetben használható lineáris egyenletrendszerek megoldására [5].

Módszer	\mathbf{G}	\mathbf{h}	\mathbf{K}	Név	Referencia
18	$\begin{bmatrix} \mathbf{A}^T \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{A} \mathbf{A}^T \end{bmatrix}$	$\begin{bmatrix} \mathbf{A}^T \mathbf{b} \\ \mathbf{A} \mathbf{c} \end{bmatrix}$	$\begin{bmatrix} \mathbf{0} & \mathbf{A}^{-1} \\ \mathbf{A}^{-T} & \mathbf{0} \end{bmatrix}$	bcr3	[7, 3.6.4]
19	$\begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}$	$\begin{bmatrix} \mathbf{A}^{-1} \mathbf{b} \\ \mathbf{A}^{-T} \mathbf{c} \end{bmatrix}$	$\begin{bmatrix} \mathbf{0} & \mathbf{A}^T \\ \mathbf{A} & \mathbf{0} \end{bmatrix}$	LSQR	[7, 3.8.2]

A 18-as módszer még egy változata a bikonjugált reziduális módszernek, amíg a 19-es módszer a Paige és Saunders féle LSQR algoritmus [5]. Itt is ortogonális transzformációk szükségesek ahhoz, hogy a megoldás következő becslését ki tudjuk számolni, mivel a gradiens számoláshoz szükséges az együtthatómátrix inverze [5].

6.4. Konklúzió

Megmutattuk, hogy nem kevesebb mint tizenkilenc konjugált gradiens algoritmust származtathatunk azzal, hogy \mathbf{G} és \mathbf{K} két szimmetrikus mátrixot választunk, és a generátorokat két különböző mód egyikével definiáljuk. Az egyik választás két-tagú (Hestenes-Stiefel) rekurrens módszerekhez vezetett, a másik 3-tagú (Lanczos) módszerekhez vezetett. A 2-tagú módszerek stabilitása garantált, ha \mathbf{G} és \mathbf{K} mátrixok definiték. Broyden a rendszertanába [5] csak 2 vagy 3-tagú rekurrens módszert vett be. Így fontos módszerek, mint pl. a GCR, GMRES, ORTHOMIN kimaradtak, és pl. az orthodir módszer csak a teljesen szimmetrikus formájában szerepelt.

7. Blokk konjugált gradiens módszer (BICG)

Jól ismert, hogy a konjugált gradiens módszereket nemszimmetrikus vagy indefinit szimmetrikus együtthatómátrixokra alkalmazva, nullával osztás miatt az algoritmus összeomlása következhet be. Ebben a fejezetben, Broyden [4] cikkére támaszkodva, szükséges és elégséges feltételeket vezetünk le az összeomlás elkerülésére, a Lanczos-féle és a Hestenes-Stiefel-féle blokk konjugált algoritmusok esetén. Ezeket a feltételeket aztán az algoritmusokat meghatározó mátrixok definittségéhez kapcsoljuk. Megmutatjuk, hogy a Lanczos-féle algoritmus stabilisabb, mint a Hestenes-Stiefel-féle algoritmus.

7.1. A Lanczos-féle és a Hestenes-Stiefel-féle algoritmusok

A fejezetben tárgyalt két változata a blokk konjugált gradiens algoritmusnak a

$$\mathbf{GX} = \mathbf{H} \quad (7.1)$$

egyenletrendszer megoldására alkalmas, ahol \mathbf{G} $n \times n$ -es valós mátrix, \mathbf{X} és \mathbf{H} pedig $n \times r$ -es valós mátrixok, ahol r általában jóval kisebb, mint n . Ha az \mathbf{X}_i mátrix az i -edik becsült megoldása a (7.1) egyenletrendszernek, és ha

$$\mathbf{F}_i = \mathbf{GX}_i - \mathbf{H}, \quad (7.2)$$

akkor

$$\mathbf{X}_{i+1} = \mathbf{X}_i + \mathbf{P}_i \mathbf{M}_i, \quad (7.3)$$

ahol

$$\mathbf{M}_i = -\mathbf{D}_i^{-1} \mathbf{P}_i^T \mathbf{F}_i \quad (7.4)$$

és

$$\mathbf{D}_i = \mathbf{P}_i^T \mathbf{G} \mathbf{P}_i. \quad (7.5)$$

Ekkor az algoritmus leállásának elégséges feltétele, a konjugáltság feltétele, azaz $\mathbf{P}_i^T \mathbf{G} \mathbf{P}_j = \mathbf{0}$, ahol $i \neq j$. Ha

$$\overline{\mathbf{P}}_i = \begin{bmatrix} \mathbf{P}_1 & \mathbf{P}_2 & \dots & \mathbf{P}_i \end{bmatrix} \quad (7.6)$$

és

$$\mathbf{Q}_i = \mathbf{I} - \mathbf{G} \overline{\mathbf{P}}_i \overline{\mathbf{D}}_i^{-1} \overline{\mathbf{P}}_i^T, \quad (7.7)$$

ahol

$$\overline{\mathbf{D}}_i = \overline{\mathbf{P}}_i^T \mathbf{G} \overline{\mathbf{P}}_i, \quad (7.8)$$

akkor $\mathbf{Q}_i \mathbf{G} \overline{\mathbf{P}_i} = \mathbf{0}$, és ha

$$\mathbf{P}_{i+1} = \mathbf{Q}_i^T \mathbf{W}_{i+1}, \quad (7.9)$$

akkor \mathbf{P}_{i+1} automatikusan kielégíti a konjugáltság feltételét bármilyen \mathbf{W}_{i+1} generátormátrixra. A $\overline{\mathbf{D}_i} = \text{diag}(\mathbf{D}_j)$ mátrix blokkdiagonális, és feltételezzük, hogy nonszinguláris.

Mivel a \mathbf{Q}_i mátrixok nagyok és ritkák, a \mathbf{P}_{i+1} mátrixok számítása a (7.9) egyenlőséggel nem praktikus. Ha a \mathbf{W}_i generátormátrixokat jól választjuk, akkor jelentős egyszerűsítések lehetségesek. A Lazncos-féle algoritmusnál

$$\mathbf{W}_i = \mathbf{K} \mathbf{G} \mathbf{P}_{i-1}, \quad (7.10)$$

ahol \mathbf{K} egy tetszőleges valós szimmetrikus mátrix, a Hestenes-Stiefel-féle algoritmusnál

$$\mathbf{W}_i = \mathbf{K} \mathbf{F}_i. \quad (7.11)$$

Ekkor a (7.9) egyenlőség a Lazncos-féle algoritmusnál

$$\mathbf{P}_{i+1} = (\mathbf{I} - \mathbf{P}_i \mathbf{D}_i^{-1} \mathbf{P}_i^T \mathbf{G} - \mathbf{P}_{i-1} \mathbf{D}_{i-1}^{-1} \mathbf{P}_{i-1}^T \mathbf{G}) \mathbf{K} \mathbf{G} \mathbf{P}_i, \quad (7.12)$$

ahol $\mathbf{P}_0 = \mathbf{0}$, \mathbf{P}_1 tetszőleges [4], és

$$\mathbf{P}_{i+1} = (\mathbf{I} - \mathbf{P}_i \mathbf{D}_i^{-1} \mathbf{P}_i^T \mathbf{G}) \mathbf{K} \mathbf{F}_{i+1} \quad (7.13)$$

a Hestenes-Stiefel-féle algoritmus esetén, ahol $\mathbf{P}_0 = \mathbf{0}$ [4].

Mindkét algoritmus stabilis, ha a \mathbf{G} és \mathbf{K} mátrixok definiték [11]. A Hestenes-Stiefel-féle algoritmus akkor és csak akkor összeomlásmentes, ha a $\mathbf{P}_i^T \mathbf{G} \mathbf{P}_i$ és az $\mathbf{F}_i^T \mathbf{K} \mathbf{F}_i$ mátrixok nonszingulárisak, ami garantált, ha a \mathbf{G} és a \mathbf{K} mátrixok definiték, és a \mathbf{P}_i és az \mathbf{F}_i mátrixok teljesrangúak [6]. Broyden a $\mathbf{P}_i^T \mathbf{G} \mathbf{P}_i$ és a $\mathbf{P}_i^T \mathbf{F}_i$ (ami $\mathbf{F} \mathbf{K} \mathbf{F}_i$ a Hestenes-Stiefel-féle algoritmus esetén) mátrixokat vizsgálta, a mögöttes Krylov sorozatok segítségével [4].

7.2. Az összeomlás elkerülésének feltételei

Legyen

$$\overline{\mathbf{V}}_i = \begin{bmatrix} \mathbf{P}_1 & \mathbf{K} \mathbf{G} \mathbf{P}_1 & (\mathbf{K} \mathbf{G})^2 \mathbf{P}_1 & \dots & (\mathbf{K} \mathbf{G})^{i-1} \mathbf{P}_1 \end{bmatrix} \quad (7.14)$$

és

$$\mathbf{S}_i = \overline{\mathbf{V}}_i^T \mathbf{G} \overline{\mathbf{V}}_i. \quad (7.15)$$

Megmutatjuk, hogy a Lanczos-féle blokk konjugált gradiens módszer akkor és csak akkor összeomlásmentes, ha az $\{\mathbf{S}_j\}$ mátrixok nonszingulárisak.

7.1. Tétel. Legyen \mathbf{P}_1 tetszőleges mátrix, és legyen $j = 1, 2, \dots, i$ -re az \mathbf{S}_j mátrix nonszinguláris, és számoljuk a \mathbf{P}_j mátrixokat a Lanczos algoritmussal (7.5-7.9 egyenlőségek, $\mathbf{W}_{i+1} = \mathbf{KGP}_i$). Ekkor $j = 1, 2, \dots, i$ -re a \mathbf{D}_j mátrixok nonszingulárisak, és a \mathbf{D}_{i+1} mátrix akkor és csak akkor nonszinguláris, ha az \mathbf{S}_{i+1} mátrix nonszinguláris. [4, 45. oldal, Theorem 1]

Bizonyítás. Lásd [4], 45. oldal. □

Tehát, ha a \mathbf{S}_j mátrixok nonszingulárisak, akkor a \mathbf{D}_j mátrixok sem azok, és így az algoritmus összemolásmentes. Másrészt, ha az \mathbf{S}_j mátrix nonszinguláris $j = 1, 2, \dots, (i-1)$ -re, akkor az \mathbf{S}_i mátrix szinguláris, ezért a \mathbf{D}_i mátrix szinguláris, és a \mathbf{P}_{i+1} és az \mathbf{X}_{i+1} mátrixok nem számolhatóak. Ezért az algoritmus összeomlik \mathbf{X}_i kiszámolása után. A Lanczos-féle algoritmusnál a összeomlásnak csak ez az egyetlen oka lehet, hogy az \mathbf{S}_i mátrix szinguláris [4].

A Hestenes-Stiefel-féle blokk konjugált gradiens módszernél a stabilitást a $\mathbf{P}_i^T \mathbf{F}_i$ mátrix is befolyásolja [6]. Legyen

$$\mathbf{T}_i = \left[\overline{\mathbf{V}}_i^T \mathbf{F}_1 \overline{\mathbf{V}}_i^T \mathbf{G} \overline{\mathbf{V}}_{i-1} \right]. \quad (7.16)$$

7.2. Tétel. A 7.1-es tétel feltételei álljanak fenn, és ezenfelül legyen az \mathbf{X}_1 mátrix tetszőleges. Ekkor, $j = 1, 2, \dots, i$ -re a $\mathbf{P}_j^T \mathbf{F}_j$ mátrix akkor és csak akkor nonszinguláris, ha a \mathbf{T}_j mátrix nonszinguláris. [4, 47. oldal, Theorem 2]

Bizonyítás. Lásd [4], 47. oldal. □

A Hestenes-Stiefel-féle blokk konjugált gradiens módszernél az algoritmus esetleges összeomlása nem csak az $\{\mathbf{S}_j\}$ mátrixsorozattól, hanem a $\{\mathbf{T}_j\}$ mátrixsorozattól is függ [4].

7.3. Tétel. Legyen \mathbf{X}_1 tetszőleges mátrix, és legyenek az \mathbf{S}_j és a \mathbf{T}_j mátrixok nonszingulárisak, $j = 1, 2, \dots, i$, és számoljuk a \mathbf{P}_{j+1} mátrixot a Hestenes-Stiefel-féle algoritmussal (7.1-7.9 egyenlőségek, $\mathbf{W}_{i+1} = \mathbf{KF}_{i+1}$). Ekkor a \mathbf{D}_j és a $\mathbf{P}_j^T \mathbf{F}_j$ mátrixok nonszingulárisak, $j = 1, 2, \dots, i$, és a \mathbf{D}_{i+1} mátrix akkor és csak akkor nonszinguláris, ha \mathbf{S}_{i+1} nonszinguláris. A $\mathbf{P}_{i+1}^T \mathbf{F}_{i+1}$ mátrix nonszinguláris, akkor és csak akkor, ha a \mathbf{T}_{i+1} mátrix nonszinguláris. [4, 48. oldal, Theorem 3]

Bizonyítás. Lásd [4], 48. oldal. □

Következmény. Ha a tétel feltételei teljesülnek, kivéve a \mathbf{T}_i mátrix nonszingularitása, akkor a \mathbf{D}_{i+1} mátrix szinguláris.

Ha az \mathbf{S}_j mátrix szinguláris, a Lanczos-féle és a Hestenes-Stiefel-féle blokk konjugált gradiens algoritmusok összeomlanak, ha a \mathbf{T}_j mátrix szinguláris, akkor csak

a Hestenes-Stiefel-féle algoritmus omlik össze. Ezért a Hestenes-Stiefel-féle blokk konjugált gradiens módszert kevésbé megbízhatónak tekinthetjük [4].

Jelölje HS felső index a Hestenes-Stiefel-féle algoritmust, és L felső index a Lanczos-féle algoritmust. Ha a kétféle módszert ugyanarra a problémára alkalmazzuk, ugyanazokkal a kezdeti értékekkel (a Lanczos-féle módszernél $\mathbf{P}_1 = \mathbf{K}\mathbf{F}_1$), akkor minden j -re $\mathbf{Q}_j^{HS} = \mathbf{Q}_j^L$ és

$$\mathbf{P}_j^{HS} = \mathbf{P}_j^L \mathbf{L}_{j-1}, \quad (7.17)$$

ahol $\mathbf{L}_{j-1} = \mathbf{M}_1^{HS} \mathbf{M}_2^{HS} \dots \mathbf{M}_{j-1}^{HS}$ [4]. A 7.3-as tétel következményéből, ha a \mathbf{T}_{j-1} mátrix szinguláris, akkor \mathbf{L}_{j-1} mátrix is szinguláris, és a (7.17) egyenlőségből következik, hogy a \mathbf{P}_j^{HS} mátrix nem teljesrangú, és így a \mathbf{D}_j^{HS} mátrix is szinguláris, akkor is, ha a $\bar{\mathbf{V}}_j$ mátrix teljesrangú. Ez egy jelentős különbség a Lanczos-féle módszerhez képest, ahol a $\bar{\mathbf{V}}_j$ és a $\bar{\mathbf{P}}_j$ mátrixok rangja mindig megegyezik [4].

Legvégül megnézzük az $\{\mathbf{X}_i\}$ mátrixsorozatot. A (7.4), (7.5) és (7.17) egyenlőségekből

$$\mathbf{M}_j^{HS} - \mathbf{L}_{j-1}^{-1} \mathbf{M}_j^L, \quad (7.18)$$

és így $\mathbf{M}_j^L = \mathbf{M}_1^{HS} \mathbf{M}_2^{HS} \dots \mathbf{M}_j^{HS}$ [4]. Ezt írhatjuk úgy is, hogy $\mathbf{M}_j^{HS} = (\mathbf{M}_{j-1}^L)^{-1} \mathbf{M}_j^L$. A (7.17) és (7.18) egyenlőségből $\mathbf{P}_j^{HS} \mathbf{M}_j^{HS} = \mathbf{P}_j^L \mathbf{M}_j^L$, így a két módszerrel generált $\{\mathbf{X}_i\}$ mátrixsorozat azonos [4].

7.3. Konklúzió

Ha $\mathbf{P}_1 = \mathbf{K}\mathbf{F}_1$, és a \mathbf{K} mátrix nemszinguláris, akkor $\mathbf{T}_i = \bar{\mathbf{V}}_i^T \mathbf{K}^{-1} \bar{\mathbf{V}}_i$ [4], és mivel $\mathbf{S}_i = \bar{\mathbf{V}}_i^T \mathbf{G} \bar{\mathbf{V}}_i$, láthatjuk, hogy [4]:

1. Ha a $\bar{\mathbf{V}}_i$ mátrix nem teljesrangú, akkor az \mathbf{S}_i és a \mathbf{T}_i mátrixok szingulárisak és mindkét algoritmus összeomlik.
2. Ha a $\bar{\mathbf{V}}_i$ mátrix teljesrangú és a \mathbf{G} mátrix indefinit, akkor mindkét algoritmusnál lehetséges az összeomlás.
3. Ha a $\bar{\mathbf{V}}_i$ mátrix teljesrangú és a \mathbf{K} mátrix indefinit, akkor a Hestenes-Stiefel-féle algoritmusnál lehetséges az összeomlás.

Az első eset kezelésére léteznek különböző módszerek [4]. A második eset komoly összeomlást eredményezhet, ha i valamilyen értékére az \mathbf{S}_i mátrix szinguláris. Ha az \mathbf{S}_i mátrix azért szinguláris, mert a $\bar{\mathbf{V}}_i$ mátrix nem teljesrangú, akkor az \mathbf{S}_k mátrix szinguláris minden $k \geq i$ esetén. Ekkor találhatunk olyan j -lépéses „előrettekintő” módszert, amivel megakadályozhatjuk az összeomlást. Ha az \mathbf{S}_i mátrix szinguláris, miközben a $\bar{\mathbf{V}}_i$ mátrix teljesrangú (mert a \mathbf{G} mátrix indefinit), akkor létezik olyan $j \geq 1$, hogy az \mathbf{S}_{i+j} mátrix nemszinguláris, vagy az \mathbf{S}_k mátrix szinguláris minden

$k \geq i$ esetén, és az összeomlás nem javítható. A harmadik esetről, ha a \mathbf{G} mátrix definit, csak a Hestenes-Stiefel-féle algoritmust használhatjuk.

Nemszimmetrikus vagy indefinit szimmetrikus együtthatómátrixok esetén sokszor megvan a lehetőségünk, hogy olyan algoritmust válasszunk, aminél vagy a \mathbf{G} , vagy a \mathbf{K} mátrix, de nem mindkét mátrix definit. A fenti elemzés alapján ilyenkor érdemes a \mathbf{K} mátrixot indefinitnek választani, és a Lanczos-féle algoritmust használni.

8. Matlab tesztfeladatokon való összehasonlítás

A Matlab egy interaktív programcsomag tudományos és mérnöki számítások, simulációk és grafikus adatmegjelenítés elvégzésére. Előre megírt függvények és eszköztárak segítik a felhasználót a gyors algoritmusfejlesztésben és az adatok grafikus megjelenítésében.

8.1. Tesztspecifikáció

A Matlab `gallery` egy tesztmátrix gyűjtemény. Az

```
[A,B,C,...] = gallery(matname,P1,P2,...)
```

függvényhívás visszaad egy tesztmátrixot vagy tesztmátrixokat, amit a `matname` paraméter specifikál. A `matname` paraméter a mátrixcsalád neve. A használható mátrixcsaládok nevét a

```
help gallery
```

paranccsal, vagy az online Matlab dokumentációban érhetjük el. A `P1,P2,...` paraméterek opcionálisak, és a `matname` által specifikált mátrixcsaládtól függenek. Ezek a paraméterek és a pontos hívási szintaktikák szintén elérhetők `help gallery` paranccsal, vagy az online Matlab dokumentációban. Ez a tesztmátrix galéria több mint ötvenféle mátrixcsaládot, típust tartalmaz, amelyek hasznosak lehetnek pl. különböző algoritmusok tesztelésénél. Ebben a fejezetben a Matlab `gallery`-vel generálunk teszt együtthatómátrixokat, és ezekkel a teszt együtthatómátrixokkal generálunk lineáris egyenletrendszereket, azaz tesztfeladatokat. A tesztfeladatok segítségével összehasonlítunk néhány, a 6. fejezetben tárgyalt konjugált gradiens módszert.

A teszt együtthatómátrixokkal úgy generálunk tesztfeladatot, hogy az \mathbf{A} teszt együtthatómátrixot egy ismert, véletlen generált \mathbf{y} vektorral megszorozzuk, és így kapunk egy \mathbf{b} vektort. A tesztfeladat az $\mathbf{Ax} = \mathbf{b}$ lineáris egyenletrendszer, ahol \mathbf{x} ismeretlen vektor. Ekkor az \mathbf{y} vektor az egyenletrendszer pontos megoldása, azaz

ismerjük a pontos eredményt, az \mathbf{x} vektor pedig mindig az egyenletrendszer adott módszerrel számolt megoldása.

A Matlab `gallery`-ből a következő mátrixcsaládokat használjuk.

- `dorr`
- `lehmer`
- `minij`
- `moler`
- `sampling`

A tesztfeladatokon a `pcg`, `SYMMLQ`, `bicgstabl`, `QMR` és az `LSQR` módszereket hasonlítjuk össze. A módszereket lásd 6. fejezet és [7]. A módszereknél a relatív reziduális toleranciáját 10^{-6} értékre állítjuk, a maximum iteráció számot pedig 1000-re. A módszereknek a saját Matlab implementációját használjuk, a függvények neve megegyezik a módszer nevével. A `pcg` és `SYMMLQ` módszereknek az $\mathbf{A}\mathbf{A}^T$ együttthatómátrixokat adtam át, mivel ezek a módszerek megkövetelik az együttthatómátrix szimmetrikusságát.

A konjugált gradiens módszereken kívül minden tesztfeladatot megoldunk a Matlab `mldivide` függvényével. Az `mldivide` függvényt az $\mathbf{x} = \mathbf{A} \backslash \mathbf{b}$ paranccsal is meg lehet hívni. Az `mldivide` megvizsgálja a feladat együttthatómátrixát, és különböző algoritmusok segítségével próbálja a lehető legjobb megoldómódszert kiválasztani. A választás mikéntje tárgyalva van a Matlab dokumentációjában.

8.2. A tesztfeladatok elemzése

A generált tesztmátrixok néhány tulajdonságát foglalja össze a következő táblázat.

Mátrixcsalád	Méret	Kondíciós szám	Rang	Pozitív definit
dorr	30×30	$5.49 * 10^5$	30	nem
lehmer	30×30	$8.66 * 10^2$	30	igen
minij	30×30	$1.50 * 10^3$	30	igen
moler	30×30	$3.18 * 10^{17}$	29	igen
sampling	30×30	$2.30 * 10^{16}$	29	nem

A mátrixcsaládokat próbáltam úgy választani, hogy a kapott tesztmátrixok különböző tulajdonságokkal rendelkezzenek.

Esetünkben két együtthatómátrix is rosszul kondicionált, vagy más szóval közel szinguláris, a moler-féle és a sampling-féle együtthatómátrixok. Az \mathbf{A} együtthatómátrix rosszul kondicionáltsága azt jelenti, hogy a \mathbf{b} vektor nagyon kicsi megváltozása az $\mathbf{Ax} = \mathbf{b}$ egyenletrendszerben, az \mathbf{x} vektor nagyon nagy megváltozását eredményezheti. A kondíciós szám „felerősíti” a pontatlanságot. Ha például a \mathbf{b} vektor elemei mérési hiba miatt ± 0.01 pontossággal ismertek, akkor az \mathbf{A} együtthatómátrix 10^2 nagyságú kondíciós száma azt jelenti, hogy az \mathbf{x} megoldásvektor elemei $\pm 100 * 0.01 = \pm 1$ pontosságúak lesznek. Természetesen ezek nem pontos definíciók. A mi tesztfeladataink mentesek bármiféle mérési hibától, mivel ismert megoldásból generáltak. A 10^{16} méretű kondíciós szám azonban már a Matlab lebegőpontos számábrázolásából eredő numerikus hibákat is felerősíti. Ez látszik abból is, hogy a Matlab `rank` függvénye szerint a moler-féle és a sampling-féle együtthatómátrixok rangja 29, pedig ez azt jelentené, hogy az együtthatómátrixok szingulárisak, és az egyenletrendszernek nincs megoldása. Valószínűleg a Matlab `rank` függvénye a rosszul kondicionáltság miatt felerősített numerikus hibák miatt hibázik. Itt még talán érdemes megjegyezni, habár a diplomamunka nem használja a determináns fogalmát, hogy egy mátrix szingularitását az alkalmazásokban nem érdemes a determináns segítségével meghatározni, mert a determináns számolásából adódó numerikus pontatlanságok miatt könnyen valótlan következtetésekre juthatunk.

8.3. Teszteredmények

A fő program végzi a tesztek, és tárolja az eredményeket. A fő program forráskódja megtalálható a 11. fejezetben.

A következő táblázatok összefoglalják a különböző módszerek eredményeit. A relatív hiba alatt az

$$\left| \frac{\|\mathbf{x}\|_2 - \|\mathbf{y}\|_2}{\|\mathbf{y}\|_2} \right|$$

értéket értjük, ahol az \mathbf{y} vektor az ismert megoldás, az \mathbf{x} vektor pedig az adott módszerrel számolt megoldás. A relatív reziduális

$$\frac{\|\mathbf{Ax}_i\|_2}{\|\mathbf{Ax}_1\|_2},$$

ahol \mathbf{x}_i az i -edik iterációhoz tartozó becslés.

A konvergencia oszlop jelentése	
0	A módszer konvergál az adott toleranciával a maximum iteráció számon belül.
1	A módszer elérte a maximum iterációt de nem konvergált.
2	A prekondicionáló mátrix rosszul kondicionált.
3	A módszer stagnált. (Két egymást követő iteráció eredménye ugyanaz)
4	Valamelyik skalár mennyiség túl kicsi, vagy túl nagy lett a folytatáshoz.

pcg	Konvergencia	Relatív reziduális	Iteráció	Relatív hiba
dorr	1	$9.57 * 10^{-3}$	1000	$8.65 * 10^{-2}$
lehmer	0	$7.60 * 10^{-7}$	30	$6.08 * 10^{-8}$
minij	0	$2.27 * 10^{-7}$	24	$4.99 * 10^{-9}$
moler	0	$3.75 * 10^{-7}$	13	$9.10 * 10^{-3}$
sampling	4	-	0	-

SYMMLQ	Konvergencia	Relatív reziduális	Iteráció	Relatív hiba
dorr	0	$2.78 * 10^{-7}$	52	$4.69 * 10^3$
lehmer	0	$4.55 * 10^{-7}$	81	$1.64 * 10^1$
minij	0	$8.85 * 10^{-7}$	50	$1.40 * 10^0$
moler	0	$1.36 * 10^{-7}$	22	$6.44 * 10^{-1}$
sampling	0	$3.03 * 10^{-7}$	18	$9.71 * 10^{-1}$

bicgstabl	Konvergencia	Relatív reziduális	Iteráció	Relatív hiba
dorr	0	$5.43 * 10^{-9}$	14	$5.43 * 10^{-5}$
lehmer	0	$6.15 * 10^{-7}$	12	$2.67 * 10^{-7}$
minij	0	$9.56 * 10^{-7}$	9	$1.13 * 10^{-5}$
moler	0	$2.50 * 10^{-7}$	5	$9.10 * 10^{-3}$
sampling	1	$2.20 * 10^{-3}$	1000	$6.64 * 10^8$

QMR	Konvergencia	Relatív reziduális	Iteráció	Relatív hiba
dorr	0	$9.63 * 10^{-11}$	30	$7.54 * 10^{-11}$
lehmer	0	$6.85 * 10^{-7}$	30	$5.38 * 10^{-8}$
minij	0	$9.94 * 10^{-7}$	22	$4.62 * 10^{-8}$
moler	0	$3.71 * 10^{-7}$	13	$9.10 * 10^{-3}$
sampling	1	$6.64 * 10^{-1}$	1000	$5.06 * 10^1$

LSQR	Konvergencia	Relatív reziduális	Iteráció	Relatív hiba
dorr	0	$8.83 * 10^{-7}$	51	$8.00 * 10^{-3}$
lehmer	0	$9.76 * 10^{-7}$	69	$1.40 * 10^{-7}$
minij	0	$7.80 * 10^{-7}$	47	$1.12 * 10^{-9}$
moler	0	$4.15 * 10^{-7}$	21	$9.10 * 10^{-3}$
sampling	0	$4.94 * 10^{-7}$	18	$1.33 * 10^{-5}$

Relatív hiba	pcg	SYMMLQ	bicgstabl	QMR	LSQR
dorr	$8.65 * 10^{-2}$	$4.69 * 10^3$	$5.43 * 10^{-5}$	$7.54 * 10^{-11}$	$8.00 * 10^{-3}$
lehmer	$6.08 * 10^{-8}$	$1.64 * 10^1$	$2.67 * 10^{-7}$	$5.38 * 10^{-8}$	$1.40 * 10^{-7}$
minij	$4.99 * 10^{-9}$	$1.40 * 10^0$	$1.13 * 10^{-5}$	$4.62 * 10^{-8}$	$1.12 * 10^{-9}$
moler	$9.10 * 10^{-3}$	$6.44 * 10^{-1}$	$9.10 * 10^{-3}$	$9.10 * 10^{-3}$	$9.10 * 10^{-3}$
sampling	-	$9.71 * 10^{-1}$	$6.64 * 10^8$	$5.06 * 10^1$	$1.33 * 10^{-5}$

	mldivide
dorr	$2.54 * 10^{-14}$
lehmer	$3.62 * 10^{-15}$
minij	$2.54 * 10^{-15}$
moler	$8.64 * 10^0$
sampling	$5.48 * 10^{-2}$

8.4. Konklúzió

A SYMMLQ módszert és a sampling-féle tesztfeladatot leszámítva összességében „jól teljesítenek” a konjugált gradiens módszerek. A sampling-féle tesztfeladat együtthatómátrixa rosszul kondicionált, kondíció száma 10^{16} , és ráadásul nem pozitív definit.

A pcg módszer a sampling-féle tesztfeladaton összeomlik, de itt az együtthatómátrix nem pozitív definit. Azonban a moler-féle tesztfeladatot $9.10 * 10^{-3}$ relatív hibával megoldja mindösszesen 13 iterációval. A moler-féle tesztfeladat együtthatómátrixa rosszul kondicionált, kondíció száma 10^{17} , de az együtthatómátrixa pozitív definit. A dorr-féle tesztfeladaton a pcg módszer nem éri el a relatív reziduális kívánt toleranciáját, de végül maximum iteráció szám mellett csak $8.65 * 10^{-2}$ lesz az eredmény relatív hibája.

A SYMMLQ módszer érdekes eredménye, hogy az első három „könnyebb” feladatot rosszabb eredménnyel oldja meg mint bármelyik másik konjugált gradiens alapú módszer, de a sampling-féle tesztfeladatot $9.71 * 10^{-1}$ relatív hibával megoldja. A moler-féle tesztfeladatot is $6.44 * 10^{-1}$ relatív hibával megoldja.

A bicgstabl módszer az első négy tesztfeladaton jól teljesít, a moler-féle tesztfeladatot $9.10 \cdot 10^{-3}$ relatív hibával megoldja, ugyanúgy mint a pcg módszer. A nem pozitív definit együtthatómátrixú sampling-féle tesztfeladaton a bicgstabl módszer nem éri el a relatív reziduális kívánt toleranciáját, és nem is tudja megoldani a tesztfeladatot, az eredmény relatív hibája $6.64 \cdot 10^8$.

A QMR módszer szintén az első négy tesztfeladatot jól megoldja. A sampling-féle tesztfeladatnál nem éri el a relatív reziduális kívánt toleranciáját, de maximum iteráció szám mellett csak $5.06 \cdot 10^1$ lesz az eredmény relatív hibája.

Az LSQR módszer kevesebb mint $8 \cdot 10^{-3}$ relatív hibával mindegyik tesztfeladatot megoldja. Az LSQR módszer pontosabban megoldja a két rosszul kondicionált együtthatómátrixos tesztfeladatot mint az `mldivide`. A sampling-féle tesztfeladatot 10^{-5} nagyságrendű relatív hibával megoldja, ami a legjobb eredmény. Ez két nagyságrenddel kisebb hiba mint az `mldivide` hibája ugyanezen a feladaton.

A mintegy referenciaként betett `mldivide` az első három tesztfeladatot messze a legkisebb hibával oldja meg. Valószínűleg csak a számábrázolás korlátozza. A relatív hibák itt 10^{-14} és 10^{-15} nagyságrendűek. A moler-féle tesztfeladatot azonban pontatlanabbul oldja meg mint bármelyik konjugált gradiens alapú módszer. A SQMMLQ feladatnál egy nagyságrenddel, a többi konjugált gradiens alapú módszerhez képest 3 nagyságrenddel pontatlanabbul. A sampling-féle tesztfeladatot $5.48 \cdot 10^{-2}$ relatív hibával oldja meg.

9. Összefoglaló

Charles George Broyden (1933. február 3. - 2011. május 20.) angol matematikus, fizikus. Jelentős szerepe volt a kvázi-Newton módszerek kifejlesztésében az 1960-as 70-es években. Az addigi módszerkehez képest jelentősen alacsonyabb számítási igényük miatt, a kvázi-Newton módszerek nagy áttörést jelentettek a nemlineáris optimalizálás területén. Broyden élete második felében a numerikus lineáris algebrára fókuszált, ezen belül is a konjugált gradiens módszerekre és ezek rendszertanára. A diplomamunka főleg a kutatásainak a konjugált gradiens módszerekkel kapcsolatos részét dolgozza fel.

Miért választottam ezt a témát? A lineáris algebra szerepe kiemelkedően fontos rengeteg területen, például jelentős geometriai, fizikai és mérnöki alkalmazásokkal rendelkezik, de a modern társadalomtudományokban is alkalmazzák. Számos más területen is találkozhatunk lineáris algebrával, szinte minden tudományág tartalmaz olyan modelleket, amelyek lineáris egyenletrendszerek megoldására vezetnek vissza valamilyen probléma megoldását.

A diplomamunka Broyden rövid életrajzával kezdődik. A diplomamunka első része a további fejezeteket alapozza meg. Ezekben a fejezetekben a további fejezetekhez

szükséges fogalmakat vezettem be, definíciókat adtam meg, tételeket mondtam ki, és tisztáztam a jelöléseket. Szintén az elméleti felvezetés része, a lineáris egyenletrendszerek megoldásának osztályozása. A direkt eljárások részben az LU-felbontást és a Gauss-féle eliminációt mutatom be részletesebben. Az iteratív eljárások utal a későbbiekre.

A további fejezetek Broyden, a numerikus lineáris algebra témakörében elért eredményeit foglalják össze.

A SOR módszer konvergencia kritériumát szimmetrikus együtthatómátrixok esetén Alexander Ostrowski dolgozta ki 1954-ben [12]. Broyden elégséges konvergencia feltételeket adott 1964-ben, szimmetrikus és nonszimmetrikus együtthatómátrixok esetére is [2]. Ezeket az eredményeket dolgozza fel a 5. fejezet.

Az egyik legjobban használható módszer az

$$\mathbf{Ax} = \mathbf{b}$$

egyenletrendszer megoldására, ahol \mathbf{A} valós nonszinguláris ritka mátrix, \mathbf{b} pedig valós vektor, a konjugált gradiens módszer és az ebből származtatott különböző módszerek [5]. Hestenes és Siefel eredeti 1952-es módszere [9] csak akkor alkalmazható, ha az \mathbf{A} együtthatómátrix szimmetrikus és pozitív definit. Az eredeti módszer óta már rengetek származtatott módszer született, amelyek nem csak szimmetrikus indefinit együtthatómátrixok esetén, de nonszimmetrikus együtthatómátrixok esetén is alkalmazhatóak [5]. A módszerek áttekintése azonban nehézkes lehet, például amiatt, hogy az algoritmusokat a különböző szerzők különböző módokon származtatják. Broyden 1996-os cikkében [5] rendszerezi a konjugált gradiens módszereket. Ezeket az eredményeket foglalja össze a 6. fejezet.

A 7. fejezetben a Lanczos-féle és a Hestenes-Stiefel-féle algortimusokat hasonlítom össze, a blokk konjugált gradiens módszeren keresztül. Broyden az összeomlás elkerüléseinek feltételeit vizsgálta [4].

A 8. fejezetben MATLAB tesztfeladatokon hasonlítom össze a pcg, SYMMLQ, bicgstabl, QMR és az LSQR módszereket. A használt Matlab kód megtalálható a függelékben.

10. Summary

Charles George Broyden (3 February 1933 – 20 May 2011) was an English mathematician, physicist. He had a major part in the development of the quasi-Newton methods in the 1960-70s. The quasi-Newton methods greatly reduced the needed computing capacity to solve nonlinear optimization problems. In the second part of his life Broyden was spending his time on numerical linear algebra. He was de-

aling mainly with the conjugate gradient methods and their taxonomy. This thesis presents his research of the second part of his life.

Why did I chose this topic? Linear algebra is really an important field. It has applications in geometry, physics, engineering, etc., it is even used in modern social sciences. We can say that we can find linear algebra in nearly every field. It is used to model real-life problems with linear systems and to obtain the solution of these systems.

The thesis begins with the short summary of Broyden's life. The first part of the thesis gives the necessary background for the following chapters. The first part of the thesis includes definitions, theorems, and clarifies the notation. It also includes a taxonomy of linear systems. The LU-decomposition and the Gauss-elimination are discussed in more detail.

The remaining chapters summarize the numerical linear algebra results of Broyden.

The convergence criteria of the SOR method for symmetric coefficient matrices were first introduced by Alexander Ostrowski in 1954 [12]. Broyden gave sufficient convergence criteria for symmetric and non-symmetric coefficient matrices in 1964 [2]. This is covered by chapter 5.

One of the best methods for solving the linear system

$$\mathbf{A}\mathbf{x} = \mathbf{b}$$

where \mathbf{A} is large and sparse is the conjugate gradient method or one of its many derivates. The original method, due to Hestenes and Stiefel [9], applies only in the case where \mathbf{A} is symmetric and positive definite but many variations have been proposed to deal not only with symmetric indefinite matrices but with non-symmetric matrices as well. However, the overview of these methods can be quite difficult because the different authors use different techniques when deriving the new algorithms. Broyden gives a new taxonomy of the conjugate gradient methods in his publication in 1996 [5]. The result is summarized in chapter 6.

In chapter 7 I compare the Lanczos and Hestenes-Stiefel algorithm through the block conjugate gradient method. Broyden investigated the breakdown conditions of these algorithms [4].

In chapter 8 I compare the pcg, SYMMLQ, bicgstabl, QMR and LSQR method in Matlab. The Matlab source code can be found at the end of the thesis.

11. Függelék

11.1. A fő program

```
number_of_testmatrices = 5;
y = randn(30, 1);

% Setup the linear systems
Test(1).Matrix = 'dorr';
Test(1).A = gallery('dorr', 30);
Test(1).b = Test(1).A*y;
Test(2).Matrix = 'lehmer';
Test(2).A = gallery('lehmer', 30);
Test(2).b = Test(2).A*y;
Test(3).Matrix = 'minij';
Test(3).A = gallery('minij', 30);
Test(3).b = Test(3).A*y;
Test(4).Matrix = 'moler';
Test(4).A = gallery('moler', 30);
Test(4).b = Test(4).A*y;
Test(5).Matrix = 'sampling';
Test(5).A = gallery('sampling', 30);
Test(5).b = Test(5).A*y;

% Save some of the properties of the coefficient matrices
for i = 1 : number_of_testmatrices
    Test(i).cond = cond(Test(i).A);
    Test(i).rank = rank(full(Test(i).A));
    if (issymmetric(Test(i).A));
        Test(i).issymmetric = 'true';
    else
        Test(i).issymmetric = 'false';
    end
    [~,p] = chol(Test(i).A);
    if (p)
        Test(i).posdef = 'false';
    else
        Test(i).posdef = 'true';
    end
end
end
```

```

for i = 1 : number_of_testmatrices
    % Calculate the result with the pcg method
    [Test(i).pcg.x, Test(i).pcg.flag, Test(i).pcg.relres, ...
    Test(i).pcg.iter] = ...
    pcg(Test(i).A, Test(i).b, 1e-6, 1000); % use A*A'
    Test(i).pcg.error = abs((norm(y) - norm(Test(i).pcg.x)) ...
    / norm(y));

    % Calculate the result with the SYMMLQ method
    [Test(i).symmlq.x, Test(i).symmlq.flag, ...
    Test(i).symmlq.relres, Test(i).symmlq.iter] = ...
    symmlq(Test(i).A*Test(i).A', Test(i).b, 1e-6, 1000);
    Test(i).symmlq.error = abs((norm(y) - norm(Test(i).symmlq.x)) ...
    / norm(y));

    % Calculate the result with the bicgstabl method
    [Test(i).bicgstabl.x, Test(i).bicgstabl.flag, ...
    Test(i).bicgstabl.relres, Test(i).bicgstabl.iter] = ...
    bicgstabl(Test(i).A, Test(i).b, 1e-6, 1000);
    Test(i).bicgstabl.error = abs((norm(y) - ...
    norm(y)) / norm(Test(i).bicgstabl.x));

    % Calculate the result with the QMR method
    [Test(i).qmr.x, Test(i).qmr.flag, Test(i).qmr.relres, ...
    Test(i).qmr.iter] = ...
    qmr(Test(i).A, Test(i).b, 1e-6, 1000);
    Test(i).qmr.error = abs((norm(y) - norm(Test(i).qmr.x)) ...
    / norm(y));

    % Calculate the result with the LSQR method
    [Test(i).lsqr.x, Test(i).lsqr.flag, Test(i).lsqr.relres, ...
    Test(i).lsqr.iter] = ...
    lsqr(Test(i).A, Test(i).b, 1e-6, 1000);
    Test(i).lsqr.error = abs((norm(y) - norm(Test(i).lsqr.x)) ...
    / norm(y));

    % Calculate the result using Matlab mldivide
    Test(i).mldivide.x = mldivide(Test(i).A, Test(i).b);

```



```

    Test(i).mldivide.error = abs((norm(y) - norm(Test(i).mldivide.x)) ...
    / norm(y));
end

```

Hivatkozások

- [1] S.F. Ashby, T.A. Manteuffel, and P.E. Saylor. A taxonomy for conjugate gradient methods. *SIAM J. Numer. Anal.*, 27:1542–1568, 1990.
- [2] C.G. Broyden. On convergence criteria for the method of successive over-relaxation. *Mathematics of Computation*, 18(85):136–141, 1964.
- [3] C.G. Broyden. *Basic Matrices*. The Macmillan Press Ltd, 1975.
- [4] C.G. Broyden. A breakdown of the block cg method. *Optimization Methods and Software*, 7(1):41–55, 1996.
- [5] C.G. Broyden. A new taxonomy of conjugate gradient methods. *Computers Math. Applic.*, 31(4/5):7–17, 1996.
- [6] C.G. Broyden and Boschetti M.A. A comparison of three basic conjugate direction methods. *Numerical Linear Algebra with Applications*, 3(6):473–489, 1996.
- [7] C.G. Broyden and M.T. Vespucchi. *Krylov Solvers for Linear Algebraic Systems*, volume 11 of *Studies in Computational Mathematics*. Elsevier B. V., 2004.
- [8] V. Faber and T. Manteuffel. Necessary and sufficient conditions for the existence of a conjugate gradient method. *SIAM J. Numer. Anal.*, 21:352–262, 1984.
- [9] M.R. Hestenes and E. Stiefel. Methods of conjugate gradients for solving linear systems. *Journal of Research of the National Bureau of Standards*, 49:409–436, 1952.
- [10] L. Nazareth. A conjugate gradient algorithm without line searches. *Journal of Optimization Theory and Applications*, 23(3):373–387, 1977.
- [11] D.P. O’Leary. The block conjugate gradient algorithm and related methods. *Linear Algebra Applic.*, 29:293–322, 1980.
- [12] A. Ostrowski. On the linear iteration procedures for symmetric matrices. *Rend. Mat. e Appl. v. 13*, page 140, 1954.
- [13] Rózsa Pál. *Lineáris algebra és alkalmazásai*. Tankönyvkiadó, Budapest, 1991.

- [14] A. Ralston. *Bevezetés a numerikus analízisbe*. Műszaki Könyvkiadó, Budapest, 1969.