

3.4 *Gyakorlat.* Az előbb kapott összefüggés segítségével bizonyítsuk be, hogy a mátrix determinánsa nem változik, ha egy oszlopának számszorosát egy másik oszlopához hozzáadjuk. Használjuk fel a szorzatmátrix determinánsára tanultakat!

3.7. Példa

Igazoljuk, hogy az $|I + ab^T|$ determináns egyenlő $1 + b^T a$ -val!

Megoldás. Feltesszük, az a és b vektorok egyike sem zérus, mert különben a feladat triviális volna. Legyen az a vektor i -edik eleme $e_i^T a = a_i \neq 0$, és tekintsük az $I - (a/a_i - e_i)e_i^T$ mátrixot. Ennek minden átlóeleme 1 és az i -edik oszlopában vannak még nemzérus elemek. De ezeket a nemzérus elemeket az i -edik sor valamely számszorosának hozzáadásával ki lehet nullázni, ebből következik, hogy a determinánsa 1. Most szorozzuk az $I + ab^T$ mátrixot balról $I - (a/a_i - e_i)e_i^T$ -vel. Ez az a vektort az $a_i e_i$ vektorba viszi, így az eredmény: $I - (a/a_i - e_i)e_i^T + a_i e_i b^T$, amely már csak az i -edik sorában és oszlopában különbözik az egységmátrixtól. Most szorozzuk a kapott mátrix k -adik oszlopát a_i/a_i -vel és adjuk hozzá a i -edik ($i \neq k$) oszlophoz (ld. 2.6 Példa):

$$\left(I - \left(\frac{a}{a_i} - e_i \right) e_i^T + a_i e_i b^T \right) \left(I + \frac{a_k}{a_i} e_k e_i^T \right) = I - \left(\frac{a - a_k e_k}{a_i} - e_i \right) e_i^T + a_i e_i b^T + a_k b_k e_i e_i^T.$$

Mint látjuk, az a vektor k -adik eleme kinullázódott, és az i -edik átlóelem $1 + a_i b_k + a_i b_k$ lett. Ezt a műveletet minden $k \neq i$ -re végrehajtva az a/a_i vektor minden átlón kívüli eleme kinullázódik, az i -edik átlóelem $1 + b^T a$, a többi pedig 1-gyel egyenlő. A következő fázisban az e_k^T , $k \neq i$ sorvektorokkal az $a_i e_i b^T$ sorvektor nemdiagonális elemeit a determináns megváltozása nélkül kinullázhatjuk.

3.8. Diádösszegek, kifejtések

Az n -edrendű egységmátrix felírható diádösszegeként: $I_n = \sum_{i=1}^n e_i e_i^T$. Ha ezt beírjuk két mátrix közé, akkor a szorzatmátrix diádösszeg-előállítását kapjuk:

$$AB = \sum_{i=1}^n A e_i e_i^T B,$$

A oszlopai és B sorai képezik a diádokat, i -edik oszlop és i -edik sor.

3.5 *Gyakorlat.* Írjuk ki ADB diádösszeg előállítását, ahol $D = [d_{ij}]$ diagonálmátrix, (csak a főátló elemei nemzérusok).

Tudjuk, az n -edrendű x vektor előállítása az egységvektorok segítségével $x = \sum_{i=1}^n e_i (e_i^T x)$. Az előállítás hasonló a $\{q_i\}_{i=1}^n$ ortonormált vektorrendszerrel. Ugyanis vezessük be a $Q = [q_1 q_2 \dots q_n]$ mátrixot. Ekkor $Q^T Q = I = Q Q^T$ az ortonormálttság miatt, tehát írható $x = Q Q^T x = \sum_{i=1}^n q_i (q_i^T x)$. Az ilyen tulajdonságú Q mátrixokat *ortogonális* (komplex megfelelője: *unitér*) mátrixoknak nevezzük.

3.9. Definíció

Az $\{a_i\}_{i=1}^n$ és $\{b_j\}_{j=1}^n$ rendszerek *biortogonális vektorrendszert* alkotnak, ha $a_i^T b_j = \alpha_i \delta_{ij}$, $\alpha_i \neq 0$ teljesül bármely indexre. Ha n a vektorok dimenziója, akkor a rendszer *teljes*.

3.6 *Gyakorlat.* Az előbbi vektorokat gyűjtjük az $A = [a_1, a_2, \dots, a_n]$ és $B = [b_1, b_2, \dots, b_n]$ mátrixba. Ellenőrizzük, hogy $A^T B$ diagonálmátrix! Ekkor az x vektor hogyan állítható elő az a_i vektorok lineáris kombinációjaként? És hogyan fejthető ki a b_j vektorok segítségével?

3.10. Tétel, mátrix egyszerű szorzatokra bontása

Minden nonszinguláris $A \in \mathbb{R}^{n \times n}$ mátrix felírható n egyszerű mátrix szorzataként, ahol egy tényező egy permutációból és egy $I + (a_i - e_i)e_i^T$ típusú tagból áll. A permutációra nincs mindig szükség.

Bizonyítás. Megadjuk az eljárást. Az első lépésben vizsgáljuk meg az A mátrix első oszlopát. Ha az első elem $a_{11} = e_1^T A e_1 \neq 0$, akkor sorcserére nincs szükség. Ha az első elem zérus, akkor az oszlopban keressünk egy nemzérus elemet, majd ennek a sorát felcseréljük az első sorral. Ha az oszlop minden eleme zérus volna, akkor nem lenne invertálható a mátrix. Az első permutáció mátrixot jelöljük Π_1 -gyel és legyen $A_1 = \Pi_1 A$.

Most szorozzuk A_1 -et a $T_1 = I - (A_1 e_1 - e_1)e_1^T / e_1^T A_1 e_1$ mátrixszal. Tudjuk, ennek eredményeként az első oszlop e_1 -be megy át és $T_1^{-1} = I + (A_1 e_1 - e_1)e_1^T$. A második lépésben hasonlóan járunk el $T_1 A_1$ második oszlopával: $A_2 = \Pi_2 T_1 A_1$ olyan mátrix lesz, ahol a 22-es pozícióban nemzérus elem van. Így a $T_2 = I - (A_2 e_2 - e_2)e_2^T / e_2^T A_2 e_2$ mátrixszal szorozva a második oszlopot az e_2 vektorba visszük. Vegyük észre, T_2 az e_1 vektort helyben hagyja.

Hasonlóan folytatva, az i -edik lépésben $A_i = \Pi_i T_{i-1} A_{i-1}$ olyan mátrix, ahol az ii pozícióban nemzérus áll. (Ha az i -edik oszlop zérus volna, ismét ellentmondásba kerülnénk azzal a feltevessel, hogy a mátrix nonszinguláris.) Ekkor a $T_i = I - (A_i e_i - e_i)e_i^T / e_i^T A_i e_i$ mátrixszal szorozva kapunk e_i -t az i -edik oszlopban és az eddig elkészült kisebb indexű egységvektorok sem romlottak el. A n -edik lépés után egységmátrixot kapunk, tehát végeredményben a mátrix inverzével szoroztunk. A szorzatokát összegyűjtve:

$$\Pi_1^T T_1^{-1} \Pi_2^T T_2^{-1} \dots T_n^{-1} = A.$$

Figyeljük meg, T_i^{-1} megadásához elég, ha az i indexet és a benne szereplő $a_i = A_i e_i$ vektort ismerjük.

3.11. Háromszögmátrixok szorzatokra bontása

Az L mátrixot alsó háromszögmátrixnak nevezzük, ha a főátló feletti elemei mind zérusok tartalmazznak. Hasonlóan az U mátrix felső háromszög mátrix, ha a főátló alatti elemek zérusok. A háromszögmátrixok szorzatokra bontása különösen egyszerű. Az előbbi tételt alkalmazva azonnal adódik az n -edrendű L alsó háromszögmátrix szorlat-előállítása:

$$L = \left(I + (L - I)e_1 e_1^T \right) \left(I + (L - I)e_2 e_2^T \right) \dots \left(I + (L - I)e_n e_n^T \right),$$

ami tömören így is írható

$$L = \prod_{i=1}^n \left(I + (L - I)e_i e_i^T \right),$$

ha megjegyezzük, hogy a tényezők növekvő indexek szerint mindig balról jobbra haladva írandók. A kifejezést közvetlenül is igazolhatjuk a j -edik oszlop meghatározásával. Jobbról az e_j vektorral szorozva az első e_j vektortól különböző eredményű szorlat $e_j + L e_j - e_j = L e_j$ az L mátrix j -edik

$$2.8. A = \begin{bmatrix} 2 & -3 & 1 \\ -4 & -2 & 1 \end{bmatrix}, \|A\|_1 = ? \quad \|A\|_\infty = ? \quad \|A\|_2 = ?$$

$$2.9. \|A\|_2 \leq \sqrt{\|A\|_1 \|A\|_\infty}.$$

2.10. Frobenius-norma: $\|A\|_F = \left(\sum_{i,j} |a_{ij}|^2 \right)^{1/2} = \sqrt{\text{tr}(A^T A)} = \sqrt{\text{tr}(A A^T)}$. Igazoljuk, ez is mátrixnorma, de nem indukált norma, $\|I\|_F = ?$, $\|Ax\|_2 \leq \|A\|_F \|x\|_2$. (A 2-es normával illeszkedő mátrixnorma.)

2.11. $A = A^T$, akkor $\|A\|_2 = \rho(A)$ = spektrál sugár, azaz szimmetrikus mátrixokra a spektrálnorma a minimális norma. ($\|\cdot\|_2$ = spektrál norma).

$$2.12. U^T U = I \text{ (ortogonális)} \rightarrow \|AU\|_F = \|A\|_F.$$

3. A numerikus lineáris algebra egyszerű transzformációi

3.1. Jelölések

A mátrixokat latin nagybetűkkel: A, B, C, \dots a vektorokat latin kisbetűkkel: a, b, c, \dots jelöljük, kivéve az i, j, k, l, m, n betűket, amelyeket indexekben fogunk használni. A skalárokat görög kisbetűket alkalmazunk. Ha az A mátrixot az a_1, a_2, \dots oszlopvektorokból állítjuk össze, akkor ezt így jelöljük: $A = [a_1 a_2 \dots a_n]$. A mátrix egy másik megadási formája $A = [a_{ij}]$, ekkor az ij -edik elemet adjuk meg általánosan. Az n -edrendű egységmátrix $I_n = [e_1 e_2 \dots e_n]$, amely az e_1, e_2, \dots, e_n Descartes-egységvektorokat tartalmazza az oszlopaiban. A transzponált jelölése: A^T , komplex esetben a transzponált konjugált jelölése A^H .

3.2. A mátrixok szorzása

Az $A = [a_{ij}] \in \mathbb{R}^{m \times n}$, $B = [b_{jk}] \in \mathbb{R}^{n \times l}$ mátrixok összeszorzásának eredménye a $C = AB = [c_{ik}] = \left[\sum_{j=1}^n a_{ij} b_{jk} \right] \in \mathbb{R}^{m \times l}$ mátrix. A vektorok egy sorból vagy oszlopból álló speciális mátrixoknak tekinthetők, szorzásuk nem jelent újat. Az alkalmazásokban megkülönböztetjük a vektorok kétféle szorzási módját. Az egyik a *skaláris* szorzat, például $a^T b$, amelynek eredménye egy skalár. A másik a *diadikus* szorzat, például ab^T , az eredmény egy 1-rangú mátrix. Vegyük észre, az első esetben szükséges, hogy a vektorok hossza azonos legyen, a második esetben nem.

3.1 Gyakorlat. Írjunk fel egy diádot. Indokoljuk meg, hogy a rangja tényleg 1. Hogyan egyszerűbb egy vektort diáddal szorozni? a) Képezzük $A = ab^T$ -t, majd Ax -et. b) Először kiszámítjuk $b^T x$ -et és ezzel a skalárral szorozzuk az a vektort.

A továbbiakban rátérünk speciális mátrixok ismertetésére.

$$(I - 2P)(I - 2P) = I - 4P + 4P = I,$$

és minden involutórus mátrix $(I - A)/2$ alakban meghatároz egy projektort:

$$(I - A)(I - A)/4 = (2I - 2A)/4 = (I - A)/2.$$

Innen látható, az egységmátrixból végtelen sokféleképp lehet gyököt vonni.

3.10 Gyakorlat. Igazoljuk, hogy a $J = [e_n, e_{n-1}, \dots, e_1]$ mátrix, ahol az egységmátrix oszlopai fordított sorrendben vannak felsorolva, involutórus mátrix. Milyen projektort határoz meg ez a mátrix, ha $n = 2, 3$?

Az $ab^T / b^T a$, $b^T a \neq 0$ projektorral a következő involutórus mátrixot készíthetjük: $I - 2ab^T / b^T a$. Az 1. ábrából megállapíthatjuk, hogy ez a mátrix a b normálisú síkra való „ferde” tükrözést végzi, ami annyit jelent, hogy az a irány mentén eljutunk a síkig, majd azt keresztezve ugyanakkora távolságot teszünk meg a túloldalon. A tükrözés akkor merőleges a síkra, ha $a = b$.

3.11 Gyakorlat. Mutassuk meg, hogy az $I - 2(x-y)(x-y)^T / (x-y)^T(x-y)$ mátrix az x és y vektorokat egymásba tükrözi, ha azok különbözőek és $x^T x = y^T y$.

3.12 Gyakorlat. Az előbbi tükröző mátrixszal lehetőségünk van arra, hogy az x vektort az $y = \pm \sigma e_1$ vektorba tükrözzük, ahol $\sigma^2 = x^T x$. Hogyan válasszuk meg σ előjelét ahhoz, hogy a kivonási jegyvesztésedet biztosan elkerüljük?

3.14. Blokk mátrixok

A mátrixokat nemcsak skalár elemekből rakhatjuk össze, hanem kisebb méretű mátrixokból is. Az ilyen mátrix elemeit *blokkoknak* nevezzük, ha pedig egy mátrixot kisebb mátrixokra bontunk, akkor a mátrixot *blokkosítjuk*. A blokkosítás történhet a következőképp: Az egységmátrixot az oszlopok mentén felszeleteljük k részre: $I = [E_1, E_2, \dots, E_k]$. Ha a sorokat ugyanilyen módon osztjuk fel blokkokra, akkor az ij -edik blokk $A_{ij} = E_i^T A E_j$ és a mátrix:

$$A = \begin{bmatrix} A_{11} & A_{12} & \dots & A_{1k} \\ A_{21} & A_{22} & \dots & A_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ A_{k1} & A_{k2} & \dots & A_{kk} \end{bmatrix}.$$

3.13 Gyakorlat. Legyen $F = I + UV^T$, U és V $n \times l$ -es mátrixok, azaz $l < n$ oszlopból állnak. Ha a kijelölt inverz létezik, ellenőrizzük: $F^{-1} = I - U(I_l + V^T U)^{-1} V^T$, ahol I_l $l \times l$ -es egységmátrix.

4. Mátrixok LU-felbontása, Gauss-Jordan algoritmus

Az LU -felbontás nem más, mint a Gauss-elimináció olyan átrendezése, ahol a részleterményeket is elrakjuk. Ez úgy történik, hogy az A mátrixot felbontjuk egy L alsó és egy U felső háromszög mátrix szorzatára.

$$\rho(A) = \max_k |\lambda_k(A)|, \quad (2.12)$$

ahol $\lambda_k(A)$ az A mátrix sajátértéke. Az $\|A\|$ mátrixnorma és az $\|x\|$ vektornorma *illeszkedő*, ha eleget tesznek a (2.8) összefüggésnek. Ez utóbbi definíció arra az esetre szól, amikor az alkalmazott mátrixnorma nem az $\|x\|$ vektornormából való indukálással készült. Igaz az összefüggés:

$$\rho(A) \leq \|A\|, \quad (2.13)$$

ha $\|A\|$ indukált vagy illeszkedő norma, mert $Au_k = \lambda_k u_k$, $\|u_k\| = 1$ mellett

$$\|A\| = \max_{\|x\|=1} \|Ax\| \geq \|Au_k\| = |\lambda_k| \|u_k\| = |\lambda_k|, \quad \forall k\text{-ra.}$$

2.9. A lineáris egyenletrendszer megoldásának perturbációi

Két esetet fogunk megvizsgálni. Az egyik, amikor az egyenletrendszer b jobboldalát perturbáljuk egy kis δb vektorral, a másik, amikor az együtthatómátrix perturbációját vizsgáljuk.

Az első esetben $A(x + \delta x) = b + \delta b$ -ből következik $A\delta x = \delta b$ és illeszkedő normák esetén kapjuk a becslést:

$$\frac{1}{\|A\| \|A^{-1}\|} \frac{\|\delta b\|}{\|b\|} \leq \frac{\|\delta x\|}{\|x\|} \leq \|A\| \|A^{-1}\| \frac{\|\delta b\|}{\|b\|}. \quad (2.14)$$

Az eredeti és a perturbált értékekre vonatkozó egyenletekből

$$\begin{array}{ccc} b = Ax, & \delta x = A^{-1} \delta b, \\ \downarrow & \downarrow \\ \|b\| \leq \|A\| \|x\|, & \|\delta x\| \leq \|A^{-1}\| \|\delta b\|. \end{array}$$

A kapott egyenlőtlenségek azonos oldalait összeszorozva kapjuk (2.14) jobboldali összefüggését. A baloldalt ugyanígy kapjuk, csak a mátrixot a másik oldalra rendezzük az induló egyenletekben:

$$\begin{array}{ccc} x = A^{-1}b, & \delta b = A\delta x, \\ \downarrow & \downarrow \\ \|x\| \leq \|A^{-1}\| \|b\|, & \|\delta b\| \leq \|A\| \|\delta x\|. \end{array}$$

Lemma. Ha $\|B\| < 1$, akkor $I + B$ invertálható és indukált normára érvényes

$$\|(I + B)^{-1}\| \leq \frac{1}{1 - \|B\|}. \quad (2.15)$$

Az előző szakaszban látott norma és spektrál sugár összefüggése szerint most B spektrál sugara kisebb 1-nél, így minden sajátértéke is kisebb, azaz nem lehet $I + B$ egyik sajátértéke sem 0.

$$(I + B)^{-1} = (I + B - B)(I + B)^{-1} = I - B(I + B)^{-1},$$

ahonnan $\|(I + B)^{-1}\| \leq 1 + \|B\| \|(I + B)^{-1}\|$, és innen átrendezéssel kapjuk az állítást. ■

Ha az együtthatómátrixot perturbáljuk δA -val: $(A + \delta A)(x + \delta x) = b \rightarrow (A + \delta A)\delta x = -\delta Ax \rightarrow \delta x = -(I + A^{-1}\delta A)^{-1} A^{-1}\delta Ax$, innen kapjuk a becslést:

$$0 \leq \frac{\|\delta x\|}{\|x\|} \leq \|(I + A^{-1}\delta A)^{-1}\| \|A^{-1}\| \|A\| \frac{\|\delta A\|}{\|A\|} \leq \|A\| \|A^{-1}\| \frac{\|\delta A\|}{\|A\|} \frac{1}{1 - \|A^{-1}\delta A\|}, \quad (2.16)$$

$$e_i^T \left(I - \begin{pmatrix} Ae_1 \\ a_{11} \end{pmatrix} e_1^T \right) Ae_k = a_{ik} - \frac{a_{11}a_{1k}}{a_{11}} = a_{ik} - \begin{pmatrix} a_{11} \\ a_{11} \end{pmatrix} a_{1k}.$$

Ez mutatja, hogy az $A - \frac{Ae_1}{e_1^T A} e_1^T A$ diádot kell számítanunk a jobb alsó $(n-1)$ -edrendű blokkra. Az

ebben szereplő oszlopvektor éppen L_1 első oszlopa, így célszerűen a következőképpen járhatunk el: kijelöljük a főelemet, vele leosztjuk az alatta lévő oszlopelemeket, a saját sorát pedig változatlanul átmásoljuk. A mátrix többi részében ebből a sorból és oszlopból készített diádot vonjuk le:

$$\begin{bmatrix} 2 & 0 & 3 & -1 \\ -4 & 5 & -2 & 3 \\ 6 & -5 & 4 & -3 \end{bmatrix} \rightarrow \begin{bmatrix} \boxed{2} & 0 & 3 & -1 \\ -2 & 5 & 4 & 1 \\ 3 & -5 & 5 & 0 \end{bmatrix} \rightarrow \begin{bmatrix} 2 & 0 & 3 & -1 \\ -2 & \boxed{5} & 4 & 1 \\ 3 & -1 & -1 & 1 \end{bmatrix}$$

$$L = \begin{bmatrix} 1 & & & \\ -2 & 1 & & \\ 3 & -1 & 1 & \end{bmatrix}, \quad U = \begin{bmatrix} 2 & 0 & 3 \\ & 5 & 4 \\ & & -1 \end{bmatrix}.$$

A végén még megoldandó $Ux = [-1 \ 1 \ 1]^T$, ezt alulról felfelé megoldva $x = [1 \ 1 \ -1]^T$.

4.1 Gyakorlat. Oldjuk meg LU-felbontással a következő egyenletrendszert:

$$\begin{bmatrix} 2 & 2 & 3 \\ 4 & 3 & 7 \\ 6 & 7 & 5 \end{bmatrix} x = \begin{bmatrix} 1 \\ 5 \\ -3 \end{bmatrix}.$$

4.2. Az LU-felbontás műveletigénye.

Az első lépésben az oszlopvektor leosztása $n-1$ osztás, a diád levonása $(n-1)^2$ szorzást és összeadást igényel. Az aritmetikai műveletek mindegyike ugyanannyi idejűnek számít, emiatt az első lépés műveletigénye: $(n-1)(2n-1)$ flop (= floating point operation, magyarul: lebegőpontos művelet). A következő lépés igénye $(n-2)(2n-3)$ flop, így a teljes műveletigény $\sum_{k=1}^{n-1} (k-1)(2k-1)$ flop. Ezt úgy közelítjük, hogy a legmagasabb fokú tagot integráljuk 0-tól n -ig: $2n^3/3$. A korrekciós tagok n kisebb hatványai, nem határozzuk meg őket, mert a legmagasabb fokú tag a domináns.

4.2 Gyakorlat. Mennyi Ax , LUx , $U^{-1}L^{-1}x$ műveletigénye? Az utolsó példánál alkalmazzuk a 2.11 szakaszban megismert faktorizációs összefüggéseket!

4.3. Blokk LU-felbontás

Néha célszerű a felbontást – vagy a mátrix invertálását – blokkosított formában végezni. Tipikusan ez a helyzet akkor, amikor az egyik elkülönített blokk egyszerűen invertálható, például azért mert egységmátrix, vagy alsó ill. felső háromszögmátrix. Mi most a blokk LU-felbontást a 2×2 -es esetben fogjuk megnézni. A főelem ilyenkor blokk, amelyről fel kell tételeznünk, hogy létezik az inverze. Legyen az egységmátrix egy felosztása $I = [E_1, E_2]$, $A_j = E_j^T A E_j$, ekkor az L mátrix a (4.1)-ben látható L_1 mátrix blokkos megfelelője (ld. még 3.13 Gyakorlat)

$$L = I - (A E_1 A_1^{-1} - E_1) E_1^T \quad (4.4)$$

és a mátrix blokkos felbontása a következő:

$$\|x\|_\infty = \max_j |x_j| \cdot \lim_{p \rightarrow \infty} \left\{ \sum_{i=1}^n \left| \frac{x_i}{\max_j |x_j|} \right|^p \right\}^{1/p} = \max_j |x_j|$$

a Csebisev-, ∞ -, vagy kocka-norma. Mint látjuk, (2.4) alapján itt a legnagyobb és legkisebb hatvány-normák szerepelnek, továbbá az ortogonális transzformációkkal szemben invariáns 2-es norma. Ezekre a normákra a definíciók alapján levezethetők a következő egyenlőtlenségek:

$$\begin{aligned} \|x\|_\infty &\leq \|x\|_1 \leq n \|x\|_\infty, \\ \|x\|_\infty &\leq \|x\|_2 \leq \sqrt{n} \|x\|_\infty, \\ \frac{1}{\sqrt{n}} \|x\|_1 &\leq \|x\|_2 \leq \|x\|_1. \end{aligned} \quad (2.5)$$

2.4. Konvergencia normában. A normák ekvivalenciája

A norma alkalmas arra, hogy segítségével egy vektorsorozat konvergenciáját értelmezzük. Ezek alapján $x^{(k)} \rightarrow x$ alatt azt értjük, hogy $\exists x \in \mathbb{R}^n$, $\lim_{k \rightarrow \infty} \|x^{(k)} - x\| = 0$.

Az $\|x\|_{(1)}$ és $\|x\|_{(2)}$ normákat *ekvivalensnek* nevezzük, ha $\exists c_1, c_2 > 0$ úgy, hogy

$$c_1 \|x\|_{(1)} \leq \|x\|_{(2)} \leq c_2 \|x\|_{(1)}.$$

6.5.1 Tétel (bizonyítás nélkül): Végesdimenziós vektortérben bármely két norma ekvivalens. Ez azt jelenti, hogy a normák akármennyire nem különbözhetnek egymástól. Így mindegy, milyen normában vizsgáljuk a konvergenciát.

2.5. Mátrixnormák

A mátrix normája $\|A\|: \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$ a következő tulajdonságokkal rendelkezik:

$$\begin{aligned} i) \quad & \|A\| = 0 \Leftrightarrow A = 0, \\ ii) \quad & \|\lambda A\| = |\lambda| \|A\|, \\ iii) \quad & \|A + B\| \leq \|A\| + \|B\|, \\ iv) \quad & \|AB\| \leq \|A\| \|B\|. \end{aligned} \quad (2.6)$$

Mivel a vektorok speciális mátrixoknak tekinthetők, minden mátrixnorma meghatároz egy vektor-normát, amelyet a mátrixnormával *kompatibilis* vektornormának nevezünk. Ez az út fordítva is bejárható, ugyanis minden vektornorma *indukál egy mátrixnormát* a következőképpen:

$$\|A\| = \sup_{x \neq 0} \frac{\|Ax\|}{\|x\|} = \sup_{\|x\|=1} \|Ax\|, \quad (2.7)$$

ahol $\|\cdot\|$ vektornorma. Csak megjegyezzük, az általánosabb definícióban megengedhető, hogy más normák szerepeljenek a számlálóban és a nevezőben. A (2.7) definíció egyenes következménye

$$\|Ax\| \leq \|A\| \|x\|. \quad (2.8)$$

2.6. Tétel

Az indukált mátrixnorma eleget tesz a (2.6) feltételeknek.

Bizonyítás. Ad 1. $A = 0 \rightarrow \|A\| = 0$. $\|A\| = 0 \rightarrow Ax = 0 \quad \forall x \rightarrow A = 0$.

alkalmazunk. Az harmadik tag azt mutatja, hogy az i -edik sort a főelemmel kell osztani, az első két tagból származó i -edik sor ugyanis zérus.

Az elmondottakat egy példán szemléltetjük. Invertálandó a $\begin{bmatrix} 0 & 1 & 1 \\ 1 & 2 & 3 \\ 1 & 3 & 6 \end{bmatrix}$ mátrix. A kibővített mátrixban

az első lépés egy sorcsere:

$$\begin{bmatrix} 0 & 1 & 1 & 0 & 0 \\ 1 & 2 & 3 & 0 & 1 \\ 1 & 3 & 6 & 0 & 1 \end{bmatrix} \xrightarrow{\substack{1 \leftrightarrow 2 \\ \text{sorcserre}}} \begin{bmatrix} 1 & 2 & 3 & 0 & 1 \\ 0 & 1 & 1 & 0 & 0 \\ 1 & 3 & 6 & 0 & 1 \end{bmatrix} \xrightarrow{Tr1}$$

Az első transzformációs lépésben az első oszlop átmegy e_1 -be, az első sort végigosztjuk a főelemmel, a többi helyen pedig végrehajtjuk az első diád levonását:

$$\rightarrow \begin{bmatrix} 1 & 2 & 3 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 & 0 \\ 0 & 1 & 3 & 0 & -1 & 1 \end{bmatrix} \xrightarrow{Tr2} \begin{bmatrix} 1 & 0 & 1 & -2 & 1 & 0 \\ 0 & 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 2 & -1 & -1 & 1 \end{bmatrix} \xrightarrow{Tr3} \begin{bmatrix} 1 & 0 & 0 & -3/2 & 3/2 & -1/2 \\ 0 & 1 & 0 & 3/2 & 1/2 & -1/2 \\ 0 & 0 & 1 & -1/2 & -1/2 & 1/2 \end{bmatrix}.$$

Az utolsó lépésben az induló egységmátrix helyén megjelent az inverz.

A „helyben” invertáláshoz azt kell észrevennünk, hogy minden lépésben összegyűjthető egy egységmátrix a kibővített mátrixból. Ezt szükségtelen tárolni. A jobboldali 3×3 -as területen minden lépésben egy új vektor jelenik meg, a bal oldali 3×3 -as területen pedig a távozó vektor helyére egy egységvektor lép be. A „tömör” algoritmusban a jobb oldalon belépő új vektort beírjuk a bal oldalon belépő egységvektor helyére. Az i -edik egységvektor helyén a jobb oldalról származó új vektor

$$\left(I - \frac{A_i e_i - e_i e_i^T}{e_i^T A_i e_i} \right) e_i = e_i - \frac{A_i e_i - e_i e_i^T}{e_i^T A_i e_i} = \begin{cases} 1/e_i^T A_i e_i, & j = i, \\ -a_{ji}^{(i)} / a_{ii}^{(i)}, & j \neq i \end{cases}$$

Ez fog átkerülni a bal oldalon az i -edik oszlopba. Így a tömör algoritmusban a főelem helyére annak reciproka kerül, az oszlop többi eleme pedig negatív előjelet kap és osztódik a főelemmel. A levonandó diád kezelése ugyanaz, mint korábban. A bekeretezett elem jelöli ki azt a diádot (sor, oszlop), amelyből a levonandó diádot képezzük. Tehát a tömör algoritmus:

$$\begin{bmatrix} 0 & 1 & 1 \\ 1 & 2 & 3 \\ 1 & 3 & 6 \end{bmatrix} \xrightarrow{\substack{1 \leftrightarrow 2 \\ \text{sorcserre}}} \begin{bmatrix} 1 & 2 & 3 \\ 0 & 1 & 1 \\ 1 & 3 & 6 \end{bmatrix} \xrightarrow{Tr1} \begin{bmatrix} 1 & 2 & 3 \\ 0 & 1 & 1 \\ -1 & 1 & 3 \end{bmatrix} \xrightarrow{Tr2} \begin{bmatrix} 1 & -2 & 1 \\ 0 & 1 & 1 \\ -1 & -1 & 2 \end{bmatrix} \xrightarrow{Tr3} \\ \rightarrow \begin{bmatrix} 3/2 & -3/2 & -1/2 \\ 1/2 & 3/2 & -1/2 \\ -1/2 & -1/2 & 1/2 \end{bmatrix} \xrightarrow{\substack{1 \leftrightarrow 2 \\ \text{oszlopcserre}}} \begin{bmatrix} -3/2 & 3/2 & -1/2 \\ 3/2 & 1/2 & -1/2 \\ -1/2 & -1/2 & 1/2 \end{bmatrix}.$$

A kezdeti sorcsere miatt nem az eredeti, hanem a ΠA mátrixot invertáltuk, ahol Π permutáció-mátrix. Ennek az inverze $A^{-1} \Pi^T$, mert $\Pi^{-1} = \Pi^T$. Így a kapott eredményt még szoroznunk kellett jobbról Π^T -vel, ami itt az $1 \leftrightarrow 2$ oszlopcserét jelenti.

4.4 Gyakorlat. Mi a Gauss-Jordan invertáló módszer műveletigényében a domináns tag?

2. Normák, egyenlőtlenségek

Ebben a szakaszban vektorok és mátrixok között távolságfüggvényeket fogunk bevezetni.

1.1. Metrikus tér

Legyen \mathcal{X} egy halmaz, amelynek elemei közt bevezetünk egy távolságfüggvényt $\delta: (\mathcal{X} \times \mathcal{X}) \rightarrow \mathbb{R}$. Azt kívánjuk, $a, b \in \mathcal{X}$ -re rendelkezzen a következő tulajdonságokkal:

- $\delta(a, b) = \delta(b, a)$, azaz a legyen olyan távolságra b -től, mint b a -tól (szimmetria).
- $\delta(a, b) = 0 \Leftrightarrow a = b$, a távolság csak akkor legyen zérus, ha a két elem azonos.
- $\delta(a, c) \leq \delta(a, b) + \delta(b, c)$, a háromszög-egyenlőtlenség. Azt fejezi ki, hogy két pont között legrövidebb út az egyenes.

Ekkor a (δ, \mathcal{X}) párt *metrikus térnek* nevezzük. A következőkben \mathcal{X} gyanánt az \mathbb{R}^n és $\mathbb{R}^{m \times n}$ halmazok kerülnek szóba, azaz vektorok és mátrixok között fogunk távolságfüggvényeket készíteni. Ez a δ nem lehet negatív, mert *iii*-ből $0 = \delta(a, a) \leq \delta(a, b) + \delta(b, a) = 2\delta(a, b)$ következmény.

2.1. A vektorok hatványnormája

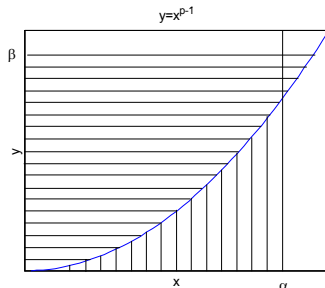
A vektor normája $\|x\|$: $\mathbb{R}^n \rightarrow \mathbb{R}$ a következő tulajdonságokkal rendelkezik:

- $\|x\| = 0 \Leftrightarrow x = 0$,
- $\|\lambda x\| = |\lambda| \|x\|$,
- $\|x + y\| \leq \|x\| + \|y\|$.

Ekkor a $\delta(x, y) = \|x - y\|$ választás metrikát ad, mert a kívánt tulajdonságok teljesülnek. Az első két feltételt triviálisan kielégíti a hatványnorma:

$$\|x\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{1/p}, \quad 1 \leq p < \infty, \quad (2.2)$$

a harmadikat később fogjuk belátni.



2.2. A Hölder-egyenlőtlenség

A hatványnormákra fennáll a Hölder-egyenlőtlenség:

$$|y^T x| \leq \sum_{i=1}^n |x_i| |y_i| \leq \|x\|_p \|y\|_q, \quad \frac{1}{p} + \frac{1}{q} = 1, \quad (2.3)$$

ami $p = q = 2$ -re a jólismert Cauchy-Bunyakovszkij egyenlőtlenségbe megy át. A p és q közötti összefüggés átrendezhető a $p - 1 = 1/(q - 1)$ alakba, amit szem előtt tartva könnyen belátható az alábbi egyenlőtlenség. Az alkalmazott függvény $y = x^{p-1}$,

az első integrál a függőlegesen, a második a vízszintesen sátriozott területet jelenti:

5.1.3 Tétel, pozitív szemidefinit mátrix felbonthatósága.

Ha A pozitív szemidefinit, akkor $A = LL^T$ alakban előállítható.

Bizonyítás. Láttuk, A főátlójában csak nemnegatív elemek lehetnek. Ha $a_{11} > 0$, akkor készítsük el a következő

$$A_2 = A - \frac{Ae_1e_1^T A}{e_1^T A e_1}, \quad (5.1)$$

mátrixot, amelyről tudjuk, hogy az első sora és oszlopa zérus. Válasszuk L első oszlopának $Le_1 = Ae_1 / \sqrt{a_{11}}$ -et, ezzel $A_2 = A - Le_1e_1^T L$.

Ha az első diagonálem zérus, akkor ugyanazon sor és oszlop cseréjével mozgassunk egy nemzérus diagonálemet az 1,1 pozícióba és ugyanígy járjunk el.

Folytassuk az eljárást a megmaradó 1-gyel kisebb méretű jobb alsó blokkal mindaddig, ameddig találunk pozitív diagonálemet. Minden lépésben az L mátrix egy újabb oszlopát nyerjük. Ha olyan helyzethez értünk, ahol a megmaradt jobb alsó blokkban minden diagonálem zérus, akkor a teljes blokknak zérusnak kell lennie, mert ha nem így volna, akkor a megmaradó blokk indefinit volna az 5.1.1 Tétel előtt tett megjegyzés szerint és ez ellentmondana annak, hogy a szemidefinités megmarad.

Vegyük észre, az alkalmazott sor-oszlop cserék a felbontást csak annyiban befolyásolják, hogy $P^T A P = LL^T$ -et kellett volna írunk, - P permutáció mátrix -, de ez átrendezhető az $A = \tilde{L}\tilde{L}^T$ alakba, ahol $\tilde{L} = PL$. ■

Szimmetrikus, pozitív definit mátrixra az $A = LL^T$ felbontást *Cholesky-felbontásnak* nevezzük. Itt most L főátlójában nem 1-esek állnak, mert például $Le_1 = Ae_1 / \sqrt{a_{11}}$ első eleme $\sqrt{a_{11}}$. A Cholesky-felbontás hasonlóképp készíthető, mint az LU -felbontás, csak most a főlemből gyököt kell vonni, és azzal végig kell osztani a saját sort és oszlopot. A számítógépes algoritmusban kihasználható, hogy a felső háromszög részt nem kell számítani, ezzel a műveletigény nagyjából megfelelődik.

5.1.4 Példa Choleski-felbontásra

$$\begin{bmatrix} 4 & -2 & 2 \\ -2 & 10 & -7 \\ 2 & -7 & 21 \end{bmatrix} \rightarrow \begin{bmatrix} \boxed{2} & -1 & 1 \\ -1 & 9 & -6 \\ 1 & -6 & 20 \end{bmatrix} \rightarrow \begin{bmatrix} 2 & & \\ -1 & \boxed{3} & -2 \\ 1 & -2 & 16 \end{bmatrix} \rightarrow \begin{bmatrix} 2 & & \\ -1 & 3 & \\ 1 & -2 & \boxed{4} \end{bmatrix},$$

$$L = \begin{bmatrix} 2 & & \\ -1 & 3 & \\ 1 & -2 & 4 \end{bmatrix}, \quad L^T = \begin{bmatrix} 2 & -1 & 1 \\ & 3 & -2 \\ & & 4 \end{bmatrix}.$$

Látható, a diádok levonása ugyanolyan módon történik, mint az LU -felbontásban.

5.1.5 Feladatok.

- Legyen $A = LL^T$ egy Cholesky felbontás. Mennyi a műveletigénye $x^T A x$ számításának, ha az eredeti mátrixot használjuk? Hogyan csökkenthető a műveletigény, ha az $x^T L L^T x$ alakot használjuk?
- A gyökvonást elkerülhetjük, ha az $A = LDL^T$ alakot használjuk, ahol L egységátlójú mátrix és D diagonálmátrix. Dolgozzuk ki ennek a felbontásnak a lépéseit! Ezt a módszert indefinit esetben is alkalmazhatjuk, ha nem adódik túlságosan kicsiny elem D -be.

csökkentésére, pl. ha $\sqrt{3472} - \sqrt{3471}$ -et így számítjuk, kihasználva, hogy a gyök alatt egész számok vannak:

$$\frac{(\sqrt{3472} - \sqrt{3471})(\sqrt{3472} + \sqrt{3471})}{\sqrt{3472} + \sqrt{3471}} = \frac{1}{\sqrt{3472} + \sqrt{3471}}.$$

A másodfokú egyenlet gyökeit pedig az alábbi módon célszerű számítani:

$$x^2 - 2px + q = 0 \text{ gyökei: } x_1 = p + \text{sign}(p)\sqrt{p^2 - q}, \quad x_2 = q/x_1.$$

- Előfordulhat olyan eset, amikor a közbülső eredmény túlsordul (nagyobb mint M_x), emiatt rossz a program futása, pedig a végeredmény az ábrázolható számok közt van. Például legyen $a = 0.3265 + 60$, $b = 0.5671 + 02$ és számítandó $\sqrt{a^2 + b^2}$. Az első szám kitevője négyzetre emeléskor 120, így túlsordult számot kapunk. Ha viszont $s\sqrt{(a/s)^2 + (b/s)^2}$ -et számítjuk, ahol $s = \max(|a|, |b|)$, akkor ez nem következik be.
- Néha arra is számítani kell, hogy egy függvény nem adja olyan pontossággal vissza a helyettesítési értéket, mint amilyen pontossággal indultunk. Például tekintsük a \sin függvényt. Ha az argumentum kicsi, akkor nincs semmi baj. Ha azonban x értéke nagy, például $x = 2356$, akkor $\sin(2356)$ számításakor 2356π -vel vett osztási maradékát kell vennünk. A maradékban már csak 1 jegy lesz pontos ha a fenti aritmetikát használjuk, így az eredménynél sem remélhetünk nagyobb pontosságot.

A mutatott példák alapján megállapíthatjuk, hogy a gépi aritmetika nemkívánatos jelenségei elsősorban akkor következnek be, ha a számok között túl nagy a nagyságrendi különbség, vagy egymáshoz nagyon közeli számokat vonunk ki egymásból.

1.5. Hibák

Az igényes számításoknál arra is kíváncsiak vagyunk, hogy az eredményt milyen pontosan tudtuk előállítani. Ehhez számba kell venni a lehetséges hibafajtákat. Az első a kiindulásul használt adatok *öröklött hibája*, nevezhetjük ezt *adathibának* is. Lehet, hogy a számítás során magunk is *tévedünk*, ezt gondos ellenőrzéssel magunknak kell felfedeznünk és kijavítanunk. A *képlethiba* az alkalmazott módszerhez tartozik. A *kerekítési hibák* részben bekövetkezhetnek a kézi számítás, adatelőkészítés során, de a gépi aritmetikának is mindig van ilyen hibája. A hibaelemzés során fel kell ismernünk, melyik az a hibafajta, ami az adott feladat szempontjából lényeges. Sok olyan számítás van, amikor az adathiba, vagy a képlethiba jelenti a fő hibaforrást. Az adathibát sokszor csak tudomásul vehetjük, de a képlethibát esetleg csökkenthetjük pontosabb módszer alkalmazásával.

A hibaszámítás alapmodellje szerint a közelítő értékekkel kapott pontos számítás eredményét közelítésnek tekintjük és azt vizsgáljuk, mekkora a hibája.

Jelölések. Az x mennyiség *pontos értéke* \hat{x} , hibája Δx : $\hat{x} = x + \Delta x$, ahol Δx előjeles szám. A relatív hiba $\delta x = \Delta x / x \approx \Delta x / \hat{x}$. Itt megjegyezzük, hogy egyes szerzők a relatív hibát a pontos értékkel definiálják, tehát az itt látható második formulát használják. A mi választásunk tudomásul veszi, hogy a pontos értéket nem ismerjük. A *hibakorlát* Δ_x egy nemnegatív szám, amellyel felülről becsüljük a hiba abszolút értékét: $|\Delta x| \leq \Delta_x$. Hasonlóképp δ_x a *relatív hibakorlát*, amelyre $|\delta x| \leq \delta_x$.

1.4 Gyakorlat. Mutassuk meg, hogy a relatív hiba kétféle megadása között a különbség másodrendű: $\delta x - \Delta x / \hat{x} = (\delta x)^2 / (1 + \delta x)$.

A valóságban a Δx hibát nem ismerjük, csak annak felső korlátját. Emiatt kiindulásul annyit tudunk, hogy \hat{x} az x érték valamely Δ_x -sugarú környezetében van.

$$K = \begin{bmatrix} 1 & & & & \\ -1 & 1 & & & \\ & \ddots & \ddots & & \\ & & & -1 & 1 \end{bmatrix}, \quad K^{-1} = S = \begin{bmatrix} 1 & & & & \\ 1 & 1 & & & \\ \vdots & \vdots & \ddots & & \\ 1 & 1 & \dots & 1 & \end{bmatrix}.$$

Inverze éppen az összegzésmátrixot adja. E két mátrix segítségével egyszerűen megadható a gyakran előforduló

$$T = \begin{bmatrix} 2 & -1 & & & \\ -1 & 2 & \ddots & & \\ & \ddots & \ddots & \ddots & \\ & & & -1 & 2 \end{bmatrix} \quad (5.2)$$

mátrix inverze:

$$T^{-1} = (K + K^T)^{-1} = [K(S + S^T)K^T]^{-1} = K^{-T}(I + ee^T)^{-1}K^{-1} = K^{-T}\left(I - \frac{ee^T}{1+n}\right)K^{-1}, \quad (5.3)$$

ahol e a csupa 1-esből álló vektor. A $T^{-1}x$ vektor előállítása így $4n$ flopp műveletet igényel.

5.3.2 Főátló-domináns háromatlós mátrix

Láttuk, ebben az esetben nem kell a főelemválasztással foglalkozni az LU -felbontás során. Ha a felbontást a mutatott módon hajtjuk végre, akkor a lineáris egyenletrendszer megoldásának műveletigénye lényegében $9n$ flopp. Háromatlós esetben van azonban két módszer is, amellyel $8n$ flopp művelettel célba érünk. A következőkben ezeket ismertetjük. Az első módszert hívhatjuk gyors LU -felbontásnak. Vegyük fel a háromatlós mátrixú egyenletrendszert a következő alakban:

$$Hx = \begin{bmatrix} d_1 & c_1 & & & \\ a_1 & d_2 & \ddots & & \\ & \ddots & \ddots & c_{n-1} & \\ & & & a_{n-1} & d_n \end{bmatrix} x = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix}. \quad (5.4)$$

Az LU -felbontás első lépése csak a második sort változtatja meg:

$$[a_1/d_1 \quad d_2 - a_1c_1/d_1 \quad c_2 \quad \dots \quad 0]x = b_2 - b_1a_1/d_1.$$

Eredményül kaptunk egy 1-gyel kisebb méretű háromatlós mátrixot, amire az eljárást megismételhetjük. Tovább folytatva végül a főelemek és jobboldalak a következők lesznek:

$$\begin{aligned} d'_1 &= d_1; & d'_i &= d_i - a_{i-1}c_{i-1}/d'_{i-1}, & i &= 2, 3, \dots, n, \\ b'_1 &= b_1; & b'_i &= b_i - a_{i-1}b'_{i-1}/d'_{i-1}, & i &= 2, 3, \dots, n. \end{aligned} \quad (5.5)$$

Most a felbontás U mátrixa felső bidiagonális - kétátlós mátrix - és a megoldandó egyenletrendszer:

$$\begin{bmatrix} d_1 & c_1 & & & \\ 0 & d'_2 & \ddots & & \\ & \ddots & \ddots & c_{n-1} & \\ & & & 0 & d'_n \end{bmatrix} x = \begin{bmatrix} b_1 \\ b'_2 \\ \vdots \\ b'_n \end{bmatrix}, \quad x_n = b'_n/d'_n; \quad x_i = b'_i - c_i x_{i+1}, \quad i = n-1, n-2, \dots, 1.$$

Láttuk: az L mátrix nem is kell a megoldáshoz, másrészt (5.5) mindkét sorában szerepel a_{i-1}/d'_{i-1} , ami elegendő egyszerű előállítani. Ezzel a megoldási algoritmus:

Például legyen a gépi számok halmaza $M(5, -4, 3)$. Ekkor a legnagyobb mantissza: $.11111 = 1 - 2^{-5}$, a legkisebb mantissza $\frac{1}{2}$. Az első pozitív gépi szám: $\varepsilon_0 = 1/2 \cdot 2^{-4} = 2^{-5}$. Az 1 után következő első gépi szám távolsága 1-től: $\varepsilon_1 = 2^{-t+1} = 2^{-4}$. A legnagyobb ábrázolható szám: $M_\infty = (1 - 2^{-t}) \cdot 2^{k^+} = (1 - 2^{-5})2^3 = 8 - 1/4$.

1.3. Valós számok konverziója gépi számmá

A következő kérdés: a valós számokat hogyan alakítsuk át gépi számokká. Az ezt megvalósító input függvényt fl-lel jelöljük (a *floating point number* kifejezés kezdőbetűi), fl: $\mathbb{R} \rightarrow M$. Megadása a következő:

$$\text{fl}(x) = \begin{cases} \infty, & \text{ha } |x| > M_\infty \\ 0, & \text{ha } |x| < \varepsilon_0 \\ x\text{-hez legközelebbi gépi szám, ha } \varepsilon_0 \leq |x| \leq M_\infty \end{cases}, \quad (1.2)$$

ahol az x -hez legközelebbi gépi szám a kerekítés szabályai szerint értendő.

Például alakítsuk át 10.87 -et 8-jegyű bináris számmá. Ezt célszerűen úgy tesszük, hogy az egész részt 2-vel osztjuk, és jegyezzük a maradékokat. A sorrendet megfordítva kapjuk a bináris jegyeket. A tört részt 2-vel szorozzuk. A kijövő egész részt nem szorozzuk tovább, hanem bináris jegyként megőrizzük. Az utolsó jegyet már abból meg tudjuk állapítani, hogy a tört rész kisebb-e 0.5 -nél. Ha kisebb, az adódó jegy 0, egyébként 1.

$$\begin{array}{r|l} 10 & 0 \\ 5 & 1 \\ 2 & 0 \\ 1 & 1 \end{array} \rightarrow 10_2 = 1010 \quad \begin{array}{r|l} . & 87 \\ 1 & 74 \\ 1 & 48 \\ 0 & 96 \end{array} \rightarrow 0.87_2 = .1101\dots$$

Kaptuk: $10.87_2 = 1010.1101\dots$. Ez nem kerekítéssel, hanem csonkítással kapott eredmény. A kerekített szám megállapításához még egy jegyet meg kell határozni. Ha a következő jegy 1, akkor az utolsó bináris jegyhez 1-et adunk, egyébként változatlanul hagyjuk. Esetünkben a következő (kilencedik) jegy 1, így a kerekített érték: 1010.1110 . Ha 10.87 -et az előbbi példában szereplő $M(5, -4, 3)$ halmazra kívánjuk leképezni az fl függvénnyel, akkor $\text{fl}(10.87) = \infty$, mert $M_\infty < 10.87$.

1.1 Gyakorlat. Legyen a gépi számok halmaza $M(5, -4, 4)$. Határozzuk meg a nevezetes számait! Mi lesz a következő számok leképezése a halmazba: $1/50$, 0.37 , 3.67 , 7.2 , 21.78 ?

1.2 Gyakorlat. Hogyan konvertálnánk 10.87 -et 3-as alapú számrendszerbe?

Feltesszük, hogy x -et pontosan ismerjük. Ekkor $\text{fl}(x)$ hibája a következőképp becsülhető:

$$|x - \text{fl}(x)| \leq \begin{cases} \infty, & \text{ha } |x| > M_\infty \\ \varepsilon_0, & \text{ha } |x| < \varepsilon_0 \\ \varepsilon_M |x|, & \text{ha } \varepsilon_0 \leq |x| \leq M_\infty \end{cases}, \quad (1.3)$$

ahol $\varepsilon_M = \varepsilon_1/2 = 2^{-t}$ a *gépi epsilon*, ez adja az ε_0 és M_∞ közé eső szám ábrázolásának relatív hibáját. Itt az első sornak csak jelzés értéke van. A második sor önmagáért beszél, egyedül a harmadik sor kíván némi magyarázatot. Azt fejezi ki, hogy az ábrázolt szám hibája nem nagyobb, mint a t -edik bináris jegyben elkövetett hiba. A harmadik sor átrendezése a relatív hiba korlátját adja:

$$\frac{|x - \text{fl}(x)|}{|x|} \leq \varepsilon_M. \quad (1.4)$$