

Seppo Linnainmaa

ALGORITMIN KUMULATIIVINEN PYÖRISTYSVIRHE

YKSITTÄISTEN PYÖRISTYSVIRHEIDEN TAYLOR-KEHITELMÄNÄ

Pro gradu-tutkielma • ohjaaja professori M.Tienari

Tapiola 1970

YHTEENVETO

Algoritmien kumulatiivista pyöristysvirhettä liukuvan pilkun aritmetiikassa pyritään analysoimaan kehittämällä tämä yksittäisten pyöristysvirheiden Taylor-kehitemmäksi. Tähän tarkoitukseen esitetään sekä analyttinen että tietokoneelle soveltuva menetelmä. Yksittäiselle suhteelliselle pyöristysvirheelle konstruoidaan tilastollinen malli. Sovellutuksina käsitellään pienalgoritmia $a^2 - b^2 = (a+b) \cdot (a-b)$ sekä Horner-shemaa ja Gauss-Jordanin matriisinkääntöalgoritmia. Sovellutuksiin liittyvät ohjelmat on ajettu IBM 7094-tietokoneella.

SISÄLLYS

Yhteenveto	II
Sisäily	III
1. Johdanto	1
2. Pyöritysvirheen käsite	2
3. Suhteellisen pyöritysvirheen jakautuma	4
- Mantissaosan pyöritysvirhe	4
- Suhteellisen pyöritysvirheen jakautuma pyörityvässä aritmetiikassa	7
- Suhteellisen pyöritysvirheen jakautuma katkaisevassa aritmetiikassa	9
4. Pyöritysvirhe laskutoimitusten yhteydessä	10
- Yhteen- ja vähennyslasku	10
- Kerto- ja jakolasku	12
5. Pyöritysvirheen Taylor-kehitemä	15
6. Taylor-sarjan analyttinen kehittäminen	19
7. Taylor-sarjan kertoimien määrittäminen tietokoneohjelmalla	24
- Kertoimienlaskualgoritmi L	24
- Algoritmia L vastaava FORTRAN IV- aliohjelmaryhmä	28
- Yksikköhäiriön menetelmä	32
8. Pienalgoritmit $a^2 - c^2$:n laskemiseksi	34
9. Horner-shema	40
- Taylor-sarjan analyttinen määrittäminen	40
- Toisen asteen kertoimet	42
- Nollakohtien lähekkäisyyden vaikutus kertoimiin	46
10. Matriisin kääntö	52
Liite: Algoritmi L FORTRAN IV-ohjelmana ...	58
Käytettyjä merkintöjä	61
Viiteluettelo	63

1. JOHDANTO

Tietokoneiden mahdollistaman suuren laskentanopeuden johdosta on tullut välttämättömäksi pyrkiä selvittämään pyöristysvirheiden vaikutuksia, koska niiden merkitys kasvaa nopeasti peräkkäisten laskutoimitusten määrän kasvaessa [1].

Kysymystä lähestyttäessä on lähtökohtana tosiasia, että algoritmista tapahtuva kumulatiivinen pyöristysvirhe on peräisin kussakin erillisessä laskutoimituksessa syntyvästä yksittäisestä pyöristysvirheestä. Näiden avulla on pyritty laskemaan ylärajoja kumulatiiviselle pyöristysvirheelle. Laajoissa algoritmeissa on ollut pakko tyytyä melko karkeisiin arvioihin tälle ylärajalle [10].

Useissa tapauksissa tieto keskimääräisestä virheestä ja sen vaihtelualttiudesta olisi hyödyllisempi kuin virheen yläraja. Tässä tarkoituksessa on pyöristysvirheitä pyritty selittämään tilastollisten jakautumien avulla. Professori Henrici on mm. teoksessaan 'Elements of Numerical Analysis' [4] laajalti lähestynyt kysymystä tältä kannalta. Hän käsittelee pääasiassa kiinteän pilkun aritmetiikkaa. Seuraavassa esityksessä pyritään tekemään vastaavia huomioita liukuvan pilkun aritmetiikasta.

2. PYÖRISTYSVIRHEEN KÄSITE

Mielivaltainen nollasta eroava kymmenjärjestelmän reaalityluku voidaan esittää muodossa

$$z = \pm 0.d_{-1}d_{-2}d_{-3}\dots \cdot 10^p, \quad (2.1)$$

missä d_{-1}, d_{-2}, \dots ovat numeroita, $d_{-1} \neq 0$ ja p on kokonaisluku. Osaa

$$m = \pm 0.d_{-1}d_{-2}d_{-3}\dots \quad (2.2)$$

kutsutaan luvun z mantissaosaksi ja sillä tarkoitetaan lukua

$$m = \pm (d_{-1} \cdot 10^{-1} + d_{-2} \cdot 10^{-2} + d_{-3} \cdot 10^{-3} + \dots). \quad (2.3)$$

Osaa 10^p kutsutaan eksponenttiosaksi, 10 on kantaluku ja p on eksponentti.

Tietokoneissa kantalukuna b on luvun 10 sijasta useimmin jokin kahden potenssi (esim. $2^1 = 2$, $2^3 = 8$ tai $2^4 = 16$). Myös tällöin luku voidaan esittää vastaavassa muodossa [4], [7]

$$z = m \cdot b^p = \pm 0.d_{-1}d_{-2}d_{-3}\dots \cdot b^p, \quad d_{-1} \neq 0, \quad (2.4)$$

missä d_{-1}, d_{-2}, \dots ovat b -kantaisen lukujärjestelmän numeroita. Tällöin mantissalla m tarkoitetaan lukua

$$m = \pm (d_{-1} \cdot b^{-1} + d_{-2} \cdot b^{-2} + d_{-3} \cdot b^{-3} + \dots). \quad (2.5)$$

Ehdosta $d_{-1} \neq 0$ seuraa

$$b^{-1} \leq |m| < 1. \quad (2.6)$$

Mantissaosaan mahtuu tietokoneessa rajoitettu määrä, t kpl, numeroita, jolloin luku z koneesta

riippuen joko katkaistaan tai pyöristetään 'oikein' muotoon

$$z^* = \pm 0.d_{-1}d_{-2}\dots d_{-t} \cdot b^p \quad (2.7)$$

Tätä muotoa kutsutaan liukuvan pilkun esitykseksi ja lukua z^* liukuluvuksi.

Pyöristetyn ja tarkan luvun erotusta

$$r = z^* - z \quad (2.8)$$

kutsutaan absoluuttiseksi pyöristysvirheeksi. Koska absoluuttisen pyöristysvirheen suuruus liukulukujen yhteydessä riippuu ratkaisevasti itse luvun z suuruudesta, on usein käytännöllisempää tarkastella suhteellista pyöristysvirhettä

$$e = \frac{z^* - z}{z} \quad , \quad (2.9)$$

jolloin

$$z^* = z(1+e) \quad . \quad (2.10)$$

3. SUHTEELLISEN PYÖRISTYSVIRHEEN JAKAUTUMA

Mantissaosan pyöristysvirhe

Suhteellisen pyöristysvirheen tilastollisen jakautuman selvittämiseksi on syytä tarkastella aluksi pelkästään mantissaosaa. Merkitsemme pyöristetyn luvun z^* mantissaosaa m^* :lla, jolloin

$$\varepsilon = m^* - m \quad (3.1)$$

on mantissaosan absoluuttinen pyöristysvirhe.

Olkoon

$$u = b^{-t} \quad (3.2)$$

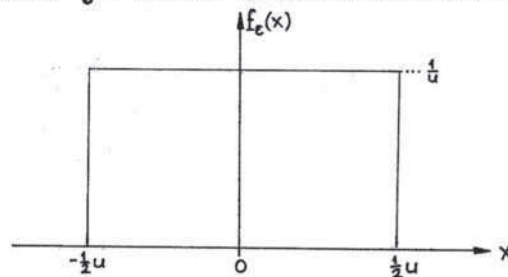
eli

$$u = 0.00\dots 01 \cdot b^0 \quad (3.3)$$

esitettyinä muodossa (2.7). Mikäli pyöristys on suoritettu 'oikein', on voimassa

$$-\frac{1}{2}u \leq \varepsilon \leq \frac{1}{2}u . \quad (3.4)$$

Intuitiivisesti pidetään selvänä, että pyöristysvirhe on satunnainen luku ja siis tilastolliselta jakautumaltaan tasaisesti jakautunut välillä $[-\frac{1}{2}u, \frac{1}{2}u]$. Satunnaisuus on tosin näennäistä: jos otamme uudelleen alkuperäisen luvun ja pyöristämme sen, saamme uudel-



kuva 1

leen saman pyöristysvirheen [4]. Tasainen jakautuma antaa kuitenkin hyvän tilastollisen mallin, jonka avulla voidaan tarkastella pyöristysvirheiden käyttäytymistä pitkissä laskutoimitussarjoissa.

Tilastollisen jakautuman tiheysfunktioilla $f(x)$ tarkoitetaan Δx :llä jaettua todennäköisyyttä sille, että muuttujan arvo on x :n ja $x+\Delta x$:n välillä, kun Δx on äärettömän pieni [2]. Pyöristysvirheen ε jakautuman tiheysfunktio on muotoa

$$f_{\varepsilon}(x) = \begin{cases} 0, & x < -\frac{1}{2}u \\ c, & -\frac{1}{2}u \leq x \leq \frac{1}{2}u \\ 0, & x > \frac{1}{2}u, \end{cases} \quad (3.5)$$

missä c on eräs vakio.

Kokonaistodennäköisyys on yksi, joten

$$\int_{-\infty}^{\infty} f_{\varepsilon}(x) dx = cu = 1, \quad (3.6)$$

mistä

$$c = \frac{1}{u}. \quad (3.7)$$

Näin saatu virheen tiheysfunktio on havainnollistettu kuvassa 1.

Virheen itseisarvolle saamme tiheysfunktioiksi

$$f_{|\varepsilon|}(x) = \begin{cases} 0, & x < 0 \\ \frac{2}{u}, & 0 \leq x \leq \frac{1}{2}u \\ 0, & x > \frac{1}{2}u. \end{cases} \quad (3.8)$$

Useissa tietokoneissa käytetään oikean pyöristyksen sijasta katkaisevaa aritmetiikkaa [6], jolloin pyöristysvirhe

$$\varepsilon' = -|m^* - m| = (m^* - m) \cdot \text{sign } m \quad (3.9)$$

on tasaisesti jakautunut välillä $[-u, 0]$. Funktio

sign x määritellään

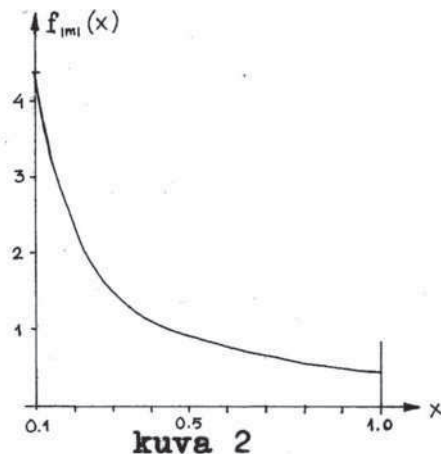
$$\text{sign } x = \begin{cases} -1, & x < 0 \\ 0, & x = 0 \\ 1, & x > 0. \end{cases} \quad (3.10)$$

Pyöristysvirheen ε' tiheysfunktiksi saamme

$$f_{\varepsilon'}(x) = \begin{cases} 0, & x < -u \\ \frac{1}{u}, & -u \leq x \leq 0 \\ 0, & x > 0. \end{cases} \quad (3.11)$$

Tarvitsemme vielä itse mantissan jakautuman pys-
tyäksemme laskemaan suhteellisen pyöristysvirheen
jakautuman. Tuntuisi luonnolliselta olettaa, että
 $|m|$ olisi tasaisesti jakautunut välillä $[b^{-1}, 1)$. On
kuitenkin todettu, ettei tämä pidä paikkaansa [3].
Jos tarkastelemme suurta joukkoa reaalimaailman kym-
menjärjestelmän lukuja, esimerkiksi fysikaalisia va-
kioita, voimme todeta, että niistä on noin 30% yk-
kösellä alkavia.

Jonkinlaisen teoreettisen selvityksen tälle ti-
lanteelle antaa samoin
ilmeisenä pitämämme to-
siasia, että kun kaikki
reaalimaailman luvut
kerrotaan tietyllä va-
kiolla, niiden ensim-
mäisten numeroiden ja-
kautuma ei muutu. Tähän
nojaamalla voidaan m :n
tiheysfunktiolle johtaa



kuvassa 2 havainnollistettu kaava [8]

$$f_{|m|}(x) = \begin{cases} 0, & x < b^{-1} \\ \frac{1}{x \cdot \ln b}, & b^{-1} \leq x < 1 \\ 0, & x \geq 1. \end{cases} \quad (3.12)$$

Tämä tiheysfunktio antaa ykkösellä alkavien lukujen esiintymistodennäköisyydeksi $\int_{0.4}^{0.2} f_{|m|}(x) dx = \log_{10} 2 = 0.30$, mikä vastaa kokeellista tulosta.

Suhteellisen pyöristysvirheen jakautuma
pyöristävässä aritmetiikassa

Kun kirjoitamme suhteellisen pyöristysvirheen e lausekkeen (2.9) mantissa- ja eksponenttiosan avulla, saamme kaavojen (2.4) ja (3.1) perusteella

$$e = \frac{z^* - z}{z} = \frac{(m^* - m) \cdot b^p}{m \cdot b^p} = \frac{\varepsilon \cdot b^p}{m \cdot b^p} = \frac{\varepsilon}{m} . \quad (3.13)$$

Johdamme nyt e :n itseisarvon $|e| = |\varepsilon/m|$ tiheysfunktion.

Kahden toisistaan riippumattoman satunnaismuuttujan η ja ξ suhteen $\zeta = \eta/\xi$ tiheysfunktio saadaan kaavasta [2]

$$f_{\zeta}(z) = \int_0^{\infty} x f_{\eta}(zx) f_{\xi}(x) dx , \quad (3.14)$$

kun $\xi > 0$. Mikäli oletamme, että ε ja m ovat riippumattomia, mikä tuntuu luonnolliselta ainakin suurilla t :n arvoilla, saadaan $|e|$:n tiheysfunktiksi kaavan (3.14) perusteella

$$f_{|e|}(z) = \begin{cases} 0 , & z < 0 \\ \int_{b^{-1}}^1 x \cdot \frac{2}{u} \cdot \frac{1}{x \cdot \ln b} dx = \frac{2 \cdot (b-1)}{ub \cdot \ln b} , & 0 \leq z \leq \frac{1}{2}u \\ \int_{b^{-1}}^{u/2z} x \cdot \frac{2}{u} \cdot \frac{1}{x \cdot \ln b} dx = \frac{1}{\ln b} \left[\frac{1}{z} - \frac{2}{ub} \right] , & \frac{1}{2}u < z \leq \frac{1}{2}ub \\ 0 , & z > \frac{1}{2}ub . \end{cases} \quad (3.15)$$

Saamme tästä $|e|$:n ylärajaksi [10]

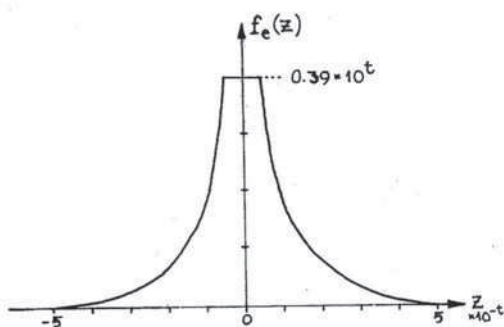
$$|e| \leq \frac{1}{2}ub . \quad (3.16)$$

Kun vielä oletamme e :n jakautuman olevan symmetrinen origon suhteen, mikä myös tuntuu ilmeiseltä,

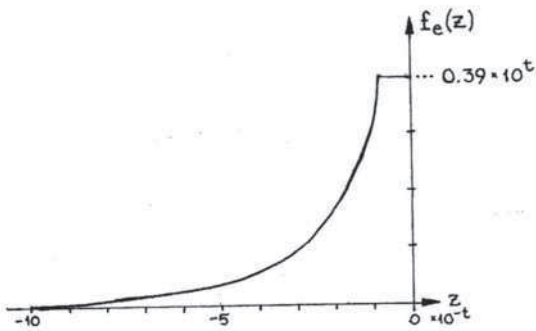
saamme e:n tiheysfunktiksi

$$f_e(z) = \begin{cases} \frac{b-1}{ub \cdot \ln b}, & |z| \leq \frac{1}{2}u \\ \frac{1}{\ln b} \left[\frac{1}{2|z|} - \frac{1}{ub} \right], & \frac{1}{2}u < |z| \leq \frac{1}{2}ub \\ 0, & |z| > \frac{1}{2}ub. \end{cases} \quad (3.17)$$

Tämä on havainnollistettu kuvassa 3a kymmenjärjestelmän tapauksessa.



kuva 3a



kuva 3b

Satunnaismuuttujan keski- eli odotusarvo saadaan kaavasta [2]

$$E(\xi) = \int_{-\infty}^{\infty} x f_{\xi}(x) dx \quad (3.18)$$

ja varianssi kaavasta

$$D^2(\xi) = \int_{-\infty}^{\infty} [x - E(\xi)]^2 f_{\xi}(x) dx. \quad (3.19)$$

Varianssin neliöjuuri, keskihajonta $D(\xi)$ kuvaa muuttujan arvojen keskimääräistä poikkeamaa keskiarvosta.

Saamme kaavojen (3.18) ja (3.19) avulla e:n odotusarvoksi ja varianssiksi

$$\mu_A = E(e) = 0, \quad \sigma_A^2 = D^2(e) = \frac{u^2}{24 \cdot \ln b} (b^2 - 1). \quad (3.20)$$

Seuraavassa taulukossa on varianssin arvoja kantalu-
vun b eri arvoilla.

b	σ_A^2/u^2
2	0.1803
4	1.262
8	1.791
10	3.832
16	41.02
32	
64	

Suhteellisen pyöristysvirheen jakautuma
katkaisevassa aritmetiikassa

Katkaisevassa aritmetiikassa saamme yhtälöstä (2.9) kaavojen (2.4) ja (3.11) perusteella suhteelliselle pyöristysvirheelle lausekkeen

$$e = \frac{z^* - z}{z} = \frac{(m^* - m) \cdot b^p}{m \cdot b^p} = \frac{m^* - m}{|m| \cdot \text{sign } m} = \frac{\varepsilon'}{|m|}. \quad (3.21)$$

Kun oletamme ε' :n ja $|m|$:n olevan riippumattomia, saamme kaavan (3.14) perusteella e :n tiheysfunktiksi (kuva 3b)

$$f_e(z) = \begin{cases} 0, & z < -ub \\ \int_{b^{-1}x}^{-u/z} \frac{1}{u} \cdot \frac{1}{x \cdot \ln b} dx = \frac{1}{\ln b} \left[-\frac{1}{z} - \frac{1}{ub} \right], & -ub \leq z < -u \\ \int_{b^{-1}x}^1 \frac{1}{u} \cdot \frac{1}{x \cdot \ln b} dx = \frac{b-1}{ub \cdot \ln b}, & -u \leq z \leq 0 \\ 0, & z > 0 \end{cases} \quad (3.22)$$

Kaavojen (3.18) ja (3.19) perusteella saamme odotusarvoksi ja varianssiksi

$$\begin{aligned} \mu_\lambda = E(e) &= \frac{u(1-b)}{2 \cdot \ln b}, \\ \sigma_\lambda^2 = D^2(e) &= -E(e) \left[\frac{u(b+1)}{3} + E(e) \right]. \end{aligned} \quad (3.23)$$

Seuraavassa taulukossa on odotusarvoja ja variansseja kantaluvun eri arvoilla.

b	μ_λ/u	σ_λ^2/u^2
2	-0.7213	0.2010
4	-1.082	
8	-1.683	2.216
10	-1.954	3.346
16	-2.705	8.011
32	-4.472	
64	-7.574	106.7

4. PYÖRISTYSVIRHE LASKUTOIMITUSTEN YHTEYDESSÄ

Yhteen- ja vähennyslasku

Kahden liukuluvun yhteenlaskun tarkkuuteen vaikuttaa ratkaisevasti yhteenlaskettavien lukujen eksponenttien p_1 ja p_2 keskinäinen ero $|p_1 - p_2|$ [10].

Mikäli $p_1 = p_2$, on yhteenlasku varsin suurella todennäköisyydellä tarkka. Ainoan poikkeuksen muodostaa tapaus, jossa saadaan muistinumero mantissojen vasemmanpuoleisimpia numeroita yhteenlaskettaessa. Tällöin joudutaan vähiten merkittävästä numerosta luopumaan, koska mantissa edelleen saa sisältää korkeintaan t numeroa. Samaan tilanteeseen saatetaan joutua myös, kun p_1 ja p_2 ovat likimain yhtäsuuret. Tämä on kuitenkin sitä epätodennäköisempää mitä suurempi $|p_1 - p_2|$ on.

Kun $|p_1 - p_2| = 1$, kohdistuu mielenkiinto lähinnä tapauksiin, joissa luvut ovat erimerkkiset ja niiden itseisarvot niin lähellä toisiaan, että tulosta joudutaan "siirtämään vasemmalle", jolloin päästään tarkkaan tulokseen.

D.W. Sweeneyn suorittamassa tutkimuksessa [8] $|p_1 - p_2|$ oli käytännön yhteenlaskuissa nolla 33-56%:n ja yksi 12-27%:n todennäköisyydellä kantaluvun saadessa arvot 2-64. Oikealle siirtojen todennäköisyys oli vastaavasti 20-2% ja vasemmalle siirtojen 20-11%.

Eksponenttien eron kasvaessa vähenee tarkan laskutoimituksen todennäköisyys ja edellä johdetut mallit pyöristävälle ja katkaisevalle aritmetiikalle pitävät varsin hyvin paikkansa. Tämä paikkansapitävyys loppuu $|p_1 - p_2|$:n ylittäessä arvon t . Tällöin luvuista itseisarvoltaan pienempi muodostaa sellaisenaan pyöristysvirheen. Näiden tapausten suhteellinen osuus kaikista yhteenlaskuista on kuitenkin niin pieni, ettei

niillä ole suhteellisen pyöristysvirheen jakautumaan merkittävää vaikutusta, kunhan huomioimme tapaukset, joissa toinen yhteenlaskettava on nolla, jolloin laskutoimitus on tarkka. Toinen operandi on edellä mainitun tutkimuksen mukaan tarkka noin 8% todennäköisyydellä.

Koska suhteellisen pyöristysvirheen jakautuma riippuu $|p_1 - p_2|:n$ jakautumasta, sille on mahdotonta johtaa täysin yleispätevää mallia. Jonkinlaisen keskiwertomallin muodostaminen tosin on mahdollista. Tällöin tulee kysymykseen lähinnä professori Tienari~~n~~ ehdottama malli: kun

$$(z_1 + z_2)^* = (z_1 + z_2) \cdot (1 + e) , \quad (4.1)$$

missä e on suhteellinen pyöristysvirhe, kirjoitetaan e muotoon

$$e = \delta \cdot e' , \quad (4.2)$$

missä e' noudattaa edellä johdettua jakautumaa ja satunnaismuuttuja δ saa arvon 0 todennäköisyydellä p sekä arvon 1 todennäköisyydellä $1-p$. Luku p edustaa tällöin todennäköisyyttä tarkalle laskutoimitukselle. Merkitsemme seuraavassa näin saatavan $e:n$ jakautuman odotusarvoa μ_s :llä ja varianssia σ_s^2 :llä.

Virheen e yläraja voidaan määrätä yleispätevästi. Perustuen epäyhtälöihin (2.6) ja (3.4) saadaan yhtälöstä (3.13) pyöristävälle aritmetiikalle [7]

$$|e| = \left| \frac{\varepsilon}{m} \right| \leq \frac{1}{2} ub. \quad (4.3)$$

Katkaisevalle aritmetiikalle saadaan vastaavasti

$$|e| \leq ub. \quad (4.4)$$

Eräissä koneissa joudutaan jo ennen laskutoimitusta pyöristämään eksponentiltaan pienempää ope-

randia. Tällöin saadaan ylärajoiksi

$$|e| \leq \frac{1}{2}u(b+1) \quad (4.5)$$

pyöristävässä ja

$$|e| \leq u(b+1) \quad (4.6)$$

katkaisevassa aritmetiikassa [10].

Kaikki yhteenlaskua koskeva koskee myös vähennyslaskua, koska näiden ero voidaan tulkita lukujen etumerkkien vaihteluksi.

Kerto- ja jakolasku

Kahden liukuluvun välinen kertolasku suoritetaan laskemalla eksponentit yhteen ja kertomalla mantissaosat keskenään. Kun

$$(z_1 \cdot z_2)^* = z_1 \cdot z_2 (1+e), \quad (4.7)$$

voidaan todeta pyöristysvirheen e noudattavan johdettua jakautumaa, so. pyöristävässä aritmetiikassa jakautumaa (3.17) ja katkaisevassa jakautumaa (3.22).

Jakolaskussa vastaavasti vähennetään eksponentit toisistaan ja mantissaosat jaetaan keskenään. Myös osamäärän pyöristysvirhe e noudattaa annettua jakautumaa, kun merkitsemme

$$\left(\frac{z_1}{z_2}\right)^* = \frac{z_1}{z_2}(1+e). \quad (4.8)$$

Merkitsemme seuraavassa kerto- ja jakolaskun pyöristysvirheen odotusarvoa μ_r :llä ja varianssia σ_r^2 :llä. Tällöin

$$\mu_r = \mu_A, \quad \sigma_r^2 = \sigma_A^2. \quad (4.9)$$

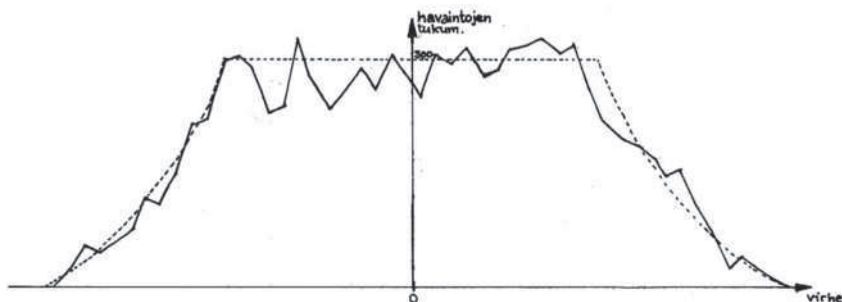
Johdetun mallin pitävyys kerto- ja jakolaskussa testattiin algoritmin S avulla.

Algoritmi S (Satunnaisalgoritmi). Algoritmi suorittaa satunnaisia laskutoimituksia $2t$ numeron tarkkuudella katkaisten tulokset t numeron mittaisiksi ja luetteloiden syntyneet pyöristysvirheet. Jos tuloksen eksponentti ylittää itseisarvoltaan puolet suurimmasta mahdollisesta eksponentista $\max |p|$, sen itseisarvoa pienennetään $\frac{1}{2} \cdot \max |p|$:llä ylivuotojen ehkäisemiseksi. Algoritmi käyttää vektoria LUKU(1), LUKU(2), ..., LUKU(100). Merkintä $A \leftarrow \text{SAT}\{B, C, \dots, Z\}$ tarkoittaa, että A:n arvoksi sijoitetaan joukon $\{B, C, \dots, Z\}$ satunnaisesti valitun alkion arvo.

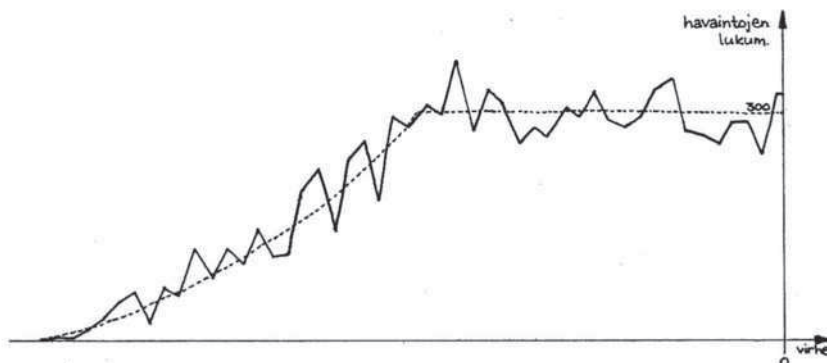
- S1. [Alkuasetukset.] Vektorin LUKU kymmenen ensimmäisen alkion arvoksi sijoitetaan tunnettuja vakioita: π , Neperin luku, Eulerin vakio, kultainen suhde jne. $N \leftarrow 10$ (N on käytössä olevan taulukon osan suurin indeksi).
- S2. [Laskutoimituksen valinta.] $I \leftarrow \text{SAT}\{1, 2, \dots, N\}$, $J \leftarrow \text{SAT}\{1, 2, \dots, N\}$, $SX \leftarrow \text{SAT}\{S3, S4, S5, S6\}$, $\rightarrow SX$.
- S3. [Yhteenlasku.] $TULOS \leftarrow \text{LUKU}(I) + \text{LUKU}(J)$, $\rightarrow S7$.
- S4. [Vähennyslasku.] $TULOS \leftarrow \text{LUKU}(I) - \text{LUKU}(J)$, $\rightarrow S7$.
- S5. [Kertolasku.] $TULOS \leftarrow \text{LUKU}(I) \cdot \text{LUKU}(J)$, $\rightarrow S7$.
- S6. [Jakolasku.] Jos $\text{LUKU}(J) = 0$, $\rightarrow S2$. Muuten $TULOS \leftarrow \text{LUKU}(I) / \text{LUKU}(J)$.
- S7. [Etumerkin valinta.] $\text{SIGN} \leftarrow \text{SAT}\{-1, 1\}$, $TULOS \leftarrow \text{SIGN} \cdot TULOS$.
- S8. [Pyöristysvirheen kirjaus.] Pyöristä mantissa t numeron mittaiseksi, luetteloi pyöristysvirheet ja pienennä tarvittaessa eksponenttia.
- S9. [Tulosalkion valinta.] Jos $N < 100$, $M \leftarrow N+1, K \leftarrow N$, muuten $K \leftarrow \text{SAT}\{1, 2, \dots, N\}$.
- S10. [Tuloksen talletus.] $\text{LUKU}(K) \leftarrow TULOS$. Jos testituloksia halutaan lisää, $\rightarrow S2$, muuten algoritmi päättyy. ■

Testitulosten lukumääräksi määrättiin 10000 sekä pyöristävälle että katkaisevalle aritmetiikalle, jolloin saatiin riittävä kuva pyöristysvirheiden jakautumasta. Summan ja erotuksen pyöristysvirheiden jakautumalle algoritmi S ei antanut mielekäästä mallia, koska eksponenttien jakautumaan ei kiinnitetty erityistä huomiota. Sen sijaan tulon ja osamäärän pyöristysvirheiden jakautumalle saatiin malli, joka vastaa teoreettista mallia. Tulokset on esitetty pyöristävälle aritmetiikalle kuvassa 4a ja katkaisevalle kuvassa 4b.

Testikoneena oli IBM 7094, jonka aritmetiikan kantalukuna on kaksi ja $t = 27$ (simuloitaessa pyöristävää aritmetiikkaa jouduttiin asettamaan $t = 26$). Mallin hyvä pitävyys osoittaa omalta osaltaan myös jakautuman (3.12) pätevyyttä kaksijärjestelmänkin yhteydessä. Tätä jakautumaahan käytettiin pyöristysvirheen teoreettista jakautumaa johdettaessa.



kuva 4a



kuva 4b

5. PYÖRISTYSVIRHEEN TAYLOR-KEHITELMÄ

Suoritettaessa peräkkäisiä laskutoimituksia liittyy jokaiseen laskutoimitukseen oma pyöristysvirheensä. Näiden yksittäisten pyöristysvirheiden yhteisvaikutusta kutsutaan kumulatiiviseksi pyöristysvirheeksi. Jos tarkastelemme esimerkiksi kahden pyöristetyn luvun a^* ja b^* jakolaskua, saamme

$$\begin{aligned}
 \left(\frac{a^*}{b^*}\right)^* &= \frac{a(1+e_a)}{b(1+e_b)}(1+e_c) \\
 &= \frac{a}{b}(1+e_a)(1+e_c)(1-e_b+e_b^2-e_b^3+\dots) \\
 &= \frac{a}{b}(1+e_a-e_b+e_c-e_a e_b+e_a e_c-e_b e_c \\
 &\quad -e_a e_b e_c+e_a e_b^2+e_c e_b^2-e_b^3+\dots) \\
 &= \frac{a}{b}(1+E_c) ,
 \end{aligned} \tag{5.1}$$

jolloin E_c edustaa kumulatiivista suhteellista pyöristysvirhettä. Kuten tässä esimerkissä voidaan kumulatiivinen pyöristysvirhe E_n aina lausua yksittäisten pyöristysvirheiden Taylor-kehittelmänä, so. muodossa

$$E_n = \sum_i a_{n,i} e_i + \sum_{i,j} a_{n,i,j} e_i e_j + \langle e^3 \rangle , \tag{5.2}$$

missä $\langle e^3 \rangle$ tarkoittaa e_i :den suhteen vähintään kolmannen asteen termejä ja a :t algoritmin alkuarvoista riippuvia vakioita. Kutsumme sarjaa (5.2) seuraavassa (E,e) -sarjaksi.

Kun merkitsemme (E,e) -sarjan k :nnetta osasummaa, so. e_i :den suhteen k :nnen asteen summaa $E_n^{(k)}$:lla, saamme yhtälön (5.2) muotoon

$$E_n = E_n^{(1)} + E_n^{(2)} + \langle e^3 \rangle . \tag{5.3}$$

Toisen ja korkeamman asteen termien vaikutus kumulatiiviseen pyöristysvirheeseen on yleensä merkityksetön yksittäisten pyöristysvirheiden e_i suuruusluokan pienuuden johdosta, joten käytännössä voidaan E_n :n lauseke ilmoittaa ensimmäisen asteen termiensä avulla [9], so.

$$E_n \approx E_n^{(1)}. \quad (5.4)$$

Kumulatiivisen pyöristysvirheen jakautuma voidaan johtaa yksittäisten pyöristysvirheiden jakautumien avulla. Jos satunnaisluvut $\xi_1, \xi_2, \dots, \xi_m$ ovat toisistaan riippumattomia ja niiden odotusarvot ovat $\mu_1, \mu_2, \dots, \mu_m$ ja varianssit $\sigma_1^2, \sigma_2^2, \dots, \sigma_m^2$, on satunnaismuuttujan

$$\xi = a_1\xi_1 + a_2\xi_2 + \dots + a_m\xi_m \quad (5.5)$$

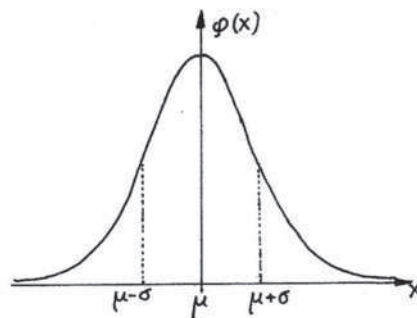
odotusarvolle ja varianssille voimassa [4]

$$E(\xi) = a_1\mu_1 + a_2\mu_2 + \dots + a_m\mu_m \quad (5.6)$$

$$D^2(\xi) \approx a_1^2\sigma_1^2 + a_2^2\sigma_2^2 + \dots + a_m^2\sigma_m^2. \quad (5.7)$$

Todennäköisyyslaskennan keskeisen raja-arvoväittämän nojalla ξ :n jakautuma lisäksi lähestyy ns. normaalijakautumaa (kuva 5) m :n lähestyessä ääretöntä [2].

Normaalijakautuman määrittävät yksikäsitteisesti sen odotusarvo ja varianssi. Sille on ominaista mm., että 68.3% havaituista satunnaisluvuista poikkeaa vähemmän kuin varianssin neliöjuuren eli keskihajonnan verran odotusarvosta.



kuva 5

Pyöristysvirheet eivät yleensä ole täysin riippumattomia toisistaan, mutta niiden väliset riippuvuudet ovat niin vähäisiä, että kaavoja (5.6) ja (5.7) käyttämällä saadaan riittäviä arvioita

kumulatiivisten pyöristysvirheiden jakautumille yksittäisten pyöristysvirheiden jakautumien avulla [4]. Oikein pyöristävässä aritmetiikassa on kaavan (5.6) perusteella myös kumulatiivisen pyöristysvirheen odotusarvo nolla.

Pitkissä laskusarjoissa voidaan nojata normaali-jakautumaan myös määritettäessä käytännöllisiä ylärajoja pyöristysvirheille. Normaalijakautumassa sijoittuu 99% kaikista havainnoista lähemmäksi odotusarvoa kuin 2.576σ , missä σ on normaalijakautuman keskihajonta. Selkeänä esimerkkinä voidaan tarkastella peräkkäisten kertolaskujen kumulatiivista pyöristysvirhettä oikein pyöristävässä aritmetiikassa.

Tuloa

$$P = \prod_{i=0}^N h_i$$

vastaa liukuvan pilkun aritmetiikkaa käyttävässä koneessa, kun alkuarvot h_i , $i = 0, \dots, N$ oletetaan tarkoiksi,

$$P^* = P(1+E_N) = \prod_{i=0}^N h_i \prod_{i=1}^N (1+e_i) \approx \prod_{i=0}^N h_i (1 + \sum_{i=1}^N e_i) ,$$

missä kukin e_i on yksittäisessä kertolaskussa syntynyt pyöristysvirhe.

Jos kertolaskut on suoritettu binääriaritmetiikassa, on virheiden e_i , $i = 1, \dots, N$ odotusarvo 0 ja varianssi $0.1803u^2$ kaavan (3.20) perusteella. Saamme E_N :n odotusarvoksi ja varianssiksi kaavojen (5.6) ja (5.7) avulla

$$E(E_N) = 0 , \quad D^2(E_N) = N \cdot 0.1803u^2 ,$$

jollöin siis 99% varmuudella

$$|E_N| < 2.576 \cdot \sqrt{N \cdot 0.1803u^2} = 1.09\sqrt{Nu} .$$

Tri Wilkinson [10] on päätenyt käytännön ylärajana tulokseen \sqrt{Nu} , mikä vastaa varsin tarkkaan edellä

johdettua tulosta. Teoreettiseksi ylärajaksi saadaan

$$|E_n| < \frac{1}{2} N \nu b ,$$

joten varsinkin suurilla $N:n$ arvoilla saataisiin varsin epärealistinen ylärajan arvo ilman tilastollista jakautumaa tai empiirisiä kokeita.

Käsitellyssä esimerkissä syntyivät kaikki pyöristysvirheet kertolaskujen yhteydessä. Pyöristysvirheet voidaan kuitenkin jakautumansa ja syntytapansa perusteella jakaa esimerkiksi tyyppeihin

1. alkuarvojen pyöristysvirheet
(keskiarvo μ_λ , varianssi σ_λ^2)
2. summan ja erotuksen pyöristysvirheet
(keskiarvo μ_s , varianssi σ_s^2)
3. tulon ja osamäärän pyöristysvirheet
(keskiarvo μ_τ , varianssi σ_τ^2).

Olkoot yhtälössä (5.2) $e_{i_1}, e_{i_2}, \dots, e_{i_k}$ tyyppiä 1, $e_{i_{k+1}}, \dots, e_{i_l}$ tyyppiä 2 ja $e_{i_{l+1}}, \dots, e_{i_m}$ tyyppiä 3 olevia pyöristysvirheitä. Tällöin $E_n:n$ odotusarvo ja varianssi voidaan kaavojen (5.4), (5.6) ja (5.7) mukaan lausua muodossa

$$E(E_n) \approx E(E_n^{(1)}) = \sum_{j=1}^k a_{n,i_j} \mu_\lambda + \sum_{j=k+1}^l a_{n,i_j} \mu_s + \sum_{j=l+1}^m a_{n,i_j} \mu_\tau , \quad (5.8)$$

$$D^2(E_n) \approx D^2(E_n^{(1)}) = \sum_{j=1}^k a_{n,i_j}^2 \sigma_\lambda^2 + \sum_{j=k+1}^l a_{n,i_j}^2 \sigma_s^2 + \sum_{j=l+1}^m a_{n,i_j}^2 \sigma_\tau^2 . \quad (5.9)$$

Esimerkiksi jakolaskun (5.1) pyöristysvirheelle E_c

$$E(E_c) \approx \mu_\tau , \quad D^2(E_c) \approx 2\sigma_\lambda^2 + \sigma_\tau^2 .$$

Mikäli alkuarvot oletetaan tarkoiksi, on yhtälössä (5.8) $\mu_\lambda = 0$ ja yhtälössä (5.9) $\sigma_\lambda^2 = 0$. Pyöristävässä aritmetiikassa $\mu_\lambda = \mu_s = \mu_\tau = 0$.

6. TAYLOR-SARJAN ANALYYTTINEN KEHITTÄMINEN

Algoritmin laskutoimitusten lukumäärän kasvaessa on yhä vaikeampaa saada selville Taylor-kehityksen kertoimia. Algoritmin rakennetta analysoimalla on kuitenkin mahdollista laskea analyyttisesti yleisiä lausekkeita kertoimille.

Kaikki algoritmit muodostuvat alkeislaskutoimituksista, joissa jokin operaatio kohdistetaan korkeintaan kahteen operandiin. Olkoon q_n erään laskutoimituksen tulos. Tällöin

$$q_n = Q_n(q_i, q_j) , \quad (6.1)$$

missä Q_n tarkoittaa yleensä yhteen-, vähennys-, kerto- tai jakolaskua. Se saattaa tarkoittaa myös potenssiin korotusta, logaritmin ottoa tms.

Tietokoneessa vastaa yhtälöä (6.1) pyöristysten vaikutuksesta yhtälö

$$q_n^* = Q_n(q_i^*, q_j^*) \cdot (1 + e_n) , \quad (6.2)$$

missä e_n on laskutoimituksen Q_n yksittäinen suhteellinen pyöristysvirhe.

Merkitsemme absoluuttista kumulatiivista pyöristysvirhettä

$$R_n = q_n^* - q_n , \quad (6.3)$$

jolloin suhteellinen kumulatiivinen pyöristysvirhe voidaan esittää muodossa

$$E_n = \frac{R_n}{q_n} . \quad (6.4)$$

Jos q_n on erityisesti alkuarvo, jolloin $q_n^* = q_n(1+e_n)$, on

$$R_n = q_n e_n \quad (6.5)$$

ja

$$E_n = e_n \quad (6.6)$$

Myös R_n voidaan lausua yksittäisten suhteellisten pyöristysvirheiden Taylor-kehitemänä

$$R_n = \sum_i c_{n,i} e_i + \sum c_{n,ij} e_i e_j + \langle e^3 \rangle \quad (6.7)$$

Nimitämme sarjaa (6.7) (R, e) -sarjaksi. Taylor-kehitemän yksikäsitteisyyden ja yhtälön (6.4) perusteella saamme E_n :n ja R_n :n yhtälöissä (5.2) ja (6.7) esiintyvälle kertoimille kaikilla kysymykseen tulevilla i :n ja j :n arvoilla

$$a_{n,i} = \frac{c_{n,i}}{q_n} \quad (6.8)$$

$$a_{n,ij} = \frac{c_{n,ij}}{q_n} \quad (6.9)$$

R_n voidaan E_n :n tapaan esittää eriasteisten osasumiensa avulla muodossa

$$R_n = R_n^{(1)} + R_n^{(2)} + \langle e^3 \rangle \quad (6.10)$$

Lausekkeen $Q_n(q_i^*, q_j^*)$ saamme väliarvolauseen ja määritelmän (6.3) nojalla muotoon [4]

$$Q_n(q_i^*, q_j^*) = q_n + D_{n,i} R_i + D_{n,j} R_j + \frac{1}{2} D_{n,ii} R_i R_i + D_{n,ij} R_i R_j + \frac{1}{2} D_{n,jj} R_j R_j + \langle R^3 \rangle, \quad (6.11)$$

missä $D_{n,i}$ on $Q_n(q_i, q_j)$:n osittaisderivaatta q_i :n suhteen pisteessä q_n ja $D_{n,ij}$ vastaavasti toinen osittaisderivaatta samassa pisteessä. Yhtälöiden (6.2),

(6.3), (6.10) ja (6.11) avulla saamme

$$R_n^{(1)} = D_{ni}R_i^{(1)} + D_{nj}R_j^{(1)} + q_n e_n \quad (6.12)$$

ja

$$R_n^{(2)} = D_{ni}R_i^{(2)} + D_{nj}R_j^{(2)} + \frac{1}{2}D_{n,ii}R_i^{(1)}R_i^{(1)} + D_{n,ij}R_i^{(1)}R_j^{(1)} + \frac{1}{2}D_{n,jj}R_j^{(1)}R_j^{(1)} + D_{ni}R_i^{(1)}e_n + D_{nj}R_j^{(1)}e_n \quad (6.13)$$

Tarkastelemme saatuja tuloksia kahden pienen esimerkin valossa. Varsinaisissa algoritmeissa päädyimme differenssiyhtälöihin, esimerkkinä tällaisista sovellamme teoriaa myöhemmin Horner-shemaan.

Esimerkki 1. Suoritetaan laskutoimitus $c = a/b$, missä a on tarkka, so. sitä ei tarvitse koneessa pyöristää, ja b :tä vastaa koneessa $b^* = b(1+e_b)$. Voimme siis yhtälön (6.5) perusteella merkitä $R_a^{(1)} = R_a^{(2)} = 0$, $R_b^{(1)} = be_b$, $R_b^{(2)} = 0$. Kaavojen (6.12) ja (6.13) perusteella

$$R_c^{(1)} = \frac{1}{b}R_a^{(1)} - \frac{a}{b^2}R_b^{(1)} + ce_c = -ce_b + ce_c$$

ja

$$R_c^{(2)} = \frac{1}{b}R_a^{(2)} - \frac{a}{b^2}R_b^{(2)} + 0 - \frac{1}{b^2}R_a^{(1)}R_b^{(1)} + \frac{a}{b^3}R_b^{(1)}R_b^{(1)} + \frac{1}{b}R_a^{(1)}e_c - \frac{a}{b^2}R_b^{(1)}e_c = ce_b e_b - ce_b e_c$$

Kaavojen (6.8), (6.9) ja (6.10) perusteella c :n (E, e) -sarja on

$$E_c = -e_b + e_c + e_b^2 - e_b e_c + \dots$$

Samaan tulokseen päädyimme aikaisemmin yhtälöissä (5.1), kun otamme huomioon, että $e_a = 0$.

Esimerkki 2. Olkoon $b^* = (\ln a^*)^*$ kun $a^* = a(1+e_a)$. Edellä esitetyt kaavat johdettiin kahdelle operandille, mutta niitä voidaan soveltaa yhden operandin operaatioihin, kunhan osittaisderivaatat 'toisen' operandin suhteen ajatellaan nolliksi, mikä on luonnollista, koska tätä ei esiinny laskutoimi-

tuksessa. Myös tätä toista operandia vastaava virhetermi ajatellaan nolllaksi. Tällöin saamme

$$R_b^{(1)} = \frac{1}{a} a e_a + \ln a \cdot e_b = e_a + b e_b$$

ja

$$R_b^{(2)} = \frac{1}{a} 0 - \frac{1}{2} \cdot \frac{1}{a^2} a^2 e_a^2 + \frac{1}{a} e_a e_b = -\frac{1}{2} e_a^2 + e_a e_b ,$$

missä e_b on logaritmin otossa tapahtuva suhteellinen virhe. Jälleen pääsemme samaan tulokseen yhtälöstä $R_b = [\ln(a+ae_a)](1+e_b) - b$, kun käytämme hyväksi kaavaa $\ln(1+e_a) = e_a - e_a^2/2 + \langle e_a^3 \rangle$.

Edellä esitetyn analyysin perusteella voimme myös verrata suhteellisten ja absoluuttisten yksittäisten pyöristysvirheiden e_n ja r_n yhteyttä kumulatiivisen pyöristysvirheen Taylor-kehityksessä.

Operaation Q_n aiheuttama yksittäinen absoluuttinen pyöristysvirhe voidaan kirjoittaa yhtälöiden (6.2) ja (6.11) mukaan muotoon

$$\begin{aligned} r_n &= Q_n(q_i^*, q_j^*) e_n \\ &= q_n e_n + D_{n,i} R_i e_n + D_{n,j} R_j e_n + \dots \end{aligned} \quad (6.14)$$

Täten R_n :n esityksestä pyöristysvirheiden r_i avulla, ns. (R,r) -sarjasta

$$R_n = \sum_i d_{n,i} r_i + \sum_{i,j} d_{n,i,j} r_i r_j + \langle r^3 \rangle \quad (6.15)$$

saadaan

$$R_n = \sum_i d_{n,i} q_i e_i + \langle e^2 \rangle . \quad (6.16)$$

Kun vertaamme yhtälöitä (6.7) ja (6.16), saamme Taylor-kehityksen yksikäsitteisyyden nojalla i :n kaikilla arvoilla

$$c_{n,i} = q_i d_{n,i} . \quad (6.17)$$

Kaavan (6.8) perusteella

$$a_{n,i} = \frac{q_i}{q_N} d_{n,i}. \quad (6.18)$$

Näin olemme selvittäneet (E,e)-, (R,e)- ja (R,r)-sarjojen ensimmäisen asteen termien väliset yhteydet.

7. TAYLOR-SARJAN KERTOIMIEN MÄÄRITTÄMINEN TIETOKONEOHJELMALLA

Kertoimienlaskualgoritmi L

Laskualgoritmien laajetessa alkaa myös Taylor-sarjan analyttinen kehittäminen tuottaa vaikeuksia, koska saadut differenssiyhtälöt monimutkaisu-
tuvat.

q_n	g_n	η_n
q_N	0	q_N
q_{N-1}	0	0
\vdots		
q_i	$D_{n,i}$	0
\vdots		
q_j	$D_{n,j}$	0
\vdots		
q_k	0	0
\vdots		
q_1	0	0

kuva 5

Voimme kuitenkin soveltaa analyttistä teoriaa tietokonealgoritmiin, jonka avulla Taylor-sarjan ensimmäisen asteen kertoimet voidaan laskea kullekin alkuarvojoukolla erikseen varsin vaikeissakin algoritmeissa.

Kertoimienlaskualgoritmi L perustuu yhtälöön (6.12). Algoritmin lähtökohtana on kaikkien alkuarvojen ja yksittäisten laskutoimitusten tulosten q_n vieminen pinon niiden esiintyessä ensimmäis-

tä kertaa. Näin syntyneen pinon (q_1, q_2, \dots, q_N) kukin elementti on siis joko alkuarvo tai tulos operaatiosta, joka on kohdistunut yhteen tai kahteen pinossa alempana olevaan elementtiin. Ajattelemmme kuhunkin elementtiin q_n , $n = 1, \dots, N$, liittyvän R_n :n ja e_n :n kertoimet g_n ja η_n siten, että ehto

$$R_N^{(a)} = \sum_{n=1}^N g_n R_n + \sum_{n=1}^N \eta_n e_n \quad (7.1)$$

on voimassa.

Aluksi toteutamme yhtälön (7.1) asettamalla $q_N = 1$, $q_n = 0$, $n \neq N$ ja $\eta_n = 0$, $n = 1, 2, \dots, N$. Olkoon $q_n = Q_n(q_i, q_j)$, missä $i > j$. Yhtälön (6.12) mukaan täytämme yhtälön (7.1) ehdon myös sijoittamalla nyt $q_N \leftarrow 0$, $q_i \leftarrow D_{n,i}$, $q_j \leftarrow D_{n,j}$ ja $\eta_n \leftarrow q_n$. Tämä tilanne on esitetty kuvassa 5.

Seuraavaksi nollaamme pinossa ylimpänä olevan nollasta poikkeavan q_n :n eli q_i :n. Jos esimerkiksi $q_i = Q_i(q_j, q_l)$, on tämä nollaus kaavan (6.12) mukaan mahdollista asettamalla $q_j \leftarrow q_j + q_i D_{i,j}$, $q_l \leftarrow q_l D_{i,l}$ ja $\eta_i \leftarrow q_i q_i$. Näin jatketaan q_n :ien nollaamista edeten pinossa alaspäin, kunnes koko pino on käyty läpi, so. $q_n = 0$, $n = 1, 2, \dots, N$. Ehto (7.1) on koko ajan voimassa, joten lopputuloksena saadaan yhtälöiden (6.7) ja (7.1) sekä Taylor-kehityksen yksikäsitteisyyden nojalla η -sarakkeeseen (R,e)-sarjan ensimmäisen asteen kertoimet.

Voimme todeta, että algoritmin edistyessä η_n pysyy nollana, kunnes q_n otetaan nollattavaksi. Tällöin η_n :n arvoksi tulee $q_n q_n$. Mikäli jätämme tässä vaiheessa kertomatta q_n :llä, saamme yhtälön (6.17) perusteella η -sarakkeeseen (R,r)-sarjan kertoimet. Voimme siis poistaa operaatiot $q_n \leftarrow 0$, $\eta_n \leftarrow q_n q_n$ ja samalla koko η -sarakkeen, jolloin algoritmin päättyessä q -sarakkeessa on (R,r)-sarjan kertoimet. Algoritmissa L vastaa q_n -kenttää kenttä COEFF(n).

Algoritmin L 1.osassa luotavan pinon I:s tietue on muotoa

VALUE(I)

TYPE	OPER1(I)	OPER2(I)
------	----------	----------

 COEFF(I) .

Kentässä VALUE on algoritmin alkuarvon tai laskutoimituksen tuloksen arvo q_I , TYPE ilmoittaa operaation Q_I laadun viereisen taulukon mukaan, OPER1 ja OPER2-kentissä on linkit operandeja vastaaviin tietueisiin ja COEFF-kent-

TYPE	operaation laatu
1	alkuarvo
2	negatointi
3	yhteenlasku
4	vähennyslasku
5	kertolasku
6	jakolasku

tään lasketaan kertoimet kuten edellä esitettiin. TYPE- ja OPER-kenttien avulla pystytään laskemaan tarvittavat osittaisderivaattojen arvot.

Kun q_i liittyy laskualgoritmissa laskutoimitukseen, siihen viitataan algoritmissa L kentän QI avulla, joka sisältää linkin q_i :n viimeksi lasketua arvoa vastaavaan tietueeseen (ko. tietueen pinoindeksiin). Seuraavassa kutsutaan QI:tä muuttujan q_i nimikkeeksi.

Algoritmi L (Kumulatiivisen pyöristysvirheen Taylor-sarjan kertoimien lasku). Algoritmi L jakautuu kahteen osaan. Mielivaltainen laskualgoritmi, 'algoritmi A', suoritetaan kokonaisuudessaan käyttäen osaa 1, joka luo tarvittavan tietuepinon. Pinolle on varattava tilaa vähintään niin monelle tietueelle kuin algoritmissa A on alkuarvoja ja laskutoimituksia yhteensä. Pinoindeksi I on pinon pohjalta lukien ensimmäisen vapaan tietueen järjestysnumero. Osassa 2 lasketaan tuloksen q_N (vastaava pinoindeksi N) kumulatiivisen pyöristysvirheen Taylor-kehityksen ensimmäisen asteen kertoimet osassa 1 muodostetun tietuepinon avulla. D1 on Q_k :n osittaisderivaatta ensimmäisen ja D2 toisen operandin suhteen voimassaolevalla pinoindeksin K arvolla. Huom. q_N :n ei välttämättä tarvitse olla algoritmin A lopputulos.

Algoritmi L, osa 1.

- L1. [Pinoindeksin alkuasetus.] $I \leftarrow 1$.
- L2. [Algoritmin A seuraava alkeistoimitus.] Jos algoritmi A on päättynyt, algoritmin L osa 1 päättyy, \rightarrow L12. Jos algoritmissa A otetaan käyttöön alkuarvo, \rightarrow L3. Jos siinä suoritetaan negointi, \rightarrow L4, jos yhteenlasku, \rightarrow L5, jos vähennyslasku, \rightarrow L6, jos kertolasku, \rightarrow L7 ja jos jakolasku, \rightarrow L8. (Tämä algoritmi ei huomioi muita toimituksia, mutta algoritmia voidaan tarvittaessa laajentaa.)

- L3. [Alkuarvo.] (Algoritmissa A otetaan käyttöön alkuarvo q_i .) $TYPE(I) \leftarrow 1$, $VALUE(I) \leftarrow q_i$, $\rightarrow L11$.
- L4. [Negatointi.] (Algoritmissa A ' $q_i = -q_j$ '.) $TYPE(I) \leftarrow 2$, $VALUE(I) \leftarrow -VALUE(QJ)$, $\rightarrow L10$.
- L5. [Yhteenlasku.] (Algoritmissa A ' $q_i = q_j + q_k$ '.) $TYPE(I) \leftarrow 3$, $VALUE(I) \leftarrow VALUE(QJ) + VALUE(QK)$, $\rightarrow L9$.
- L6. [Vähennyslasku.] (Algoritmissa A ' $q_i = q_j - q_k$ '.) $TYPE(I) \leftarrow 4$, $VALUE(I) \leftarrow VALUE(QJ) - VALUE(QK)$, $\rightarrow L9$.
- L7. [Kertolasku.] (Algoritmissa A ' $q_i = q_j \cdot q_k$ '.) $TYPE(I) \leftarrow 5$, $VALUE(I) \leftarrow VALUE(QJ) * VALUE(QK)$, $\rightarrow L9$.
- L8. [Jakolasku.] (Algoritmissa A ' $q_i = q_j / q_k$ '.) $TYPE(I) \leftarrow 6$, $VALUE(I) \leftarrow VALUE(QJ) / VALUE(QK)$. $\rightarrow L9$.
- L9. [Linkki 2. operandiin.] $OPER2(I) \leftarrow QK$.
- L10. [Linkki 1. operandiin.] $OPER1(I) \leftarrow QJ$.
- L11. [Nimike kuntoon.] $QI \leftarrow I$, $I \leftarrow I + 1$, $\rightarrow L2$.

Algoritmi L, osa 2.

- L12. [COEFF-kenttien alkuasetus.] Nollaa kentät $COEFF(K)$, $K = 1, 2, \dots, N-1$. $COEFF(N) \leftarrow 1$, $K \leftarrow N$.
- L13. [Tyypivalinta.] Mene askeleeseen LX, missä $X = 13 + TYPE(K)$.
- L14. [Alkuarvo.] $\rightarrow L21$.
- L15. [Negatointi.] $COEFF(OPER1(K)) \leftarrow -COEFF(OPER1(K)) - COEFF(K)$, $COEFF(K) \leftarrow 0$ (koska negatointi on tarkka toimitus), $\rightarrow L21$.
- L16. [Yhteenlasku.] $D1 \leftarrow 1$, $D2 \leftarrow 1$, $\rightarrow L20$.
- L17. [Vähennyslasku.] $D1 \leftarrow 1$, $D2 \leftarrow -1$, $\rightarrow L20$.
- L18. [Kertolasku.] $D1 \leftarrow VALUE(OPER2(K))$, $D2 \leftarrow VALUE(OPER1(K))$, $\rightarrow L20$.
- L19. [Jakolasku.] $D1 \leftarrow 1 / VALUE(OPER2(K))$, $D2 \leftarrow -VALUE(OPER1(K)) / VALUE(OPER2(K)) ** 2$.

- L20. [COEFF-kenttien käsittely.]
COEFF(OPER1(K))←COEFF(OPER1(K))+D1*COEFF(K),
COEFF(OPER2(K))←COEFF(OPER2(K))+D2*COEFF(K).
- L21. [Indeksin vähennys.] $K \leftarrow K-1$. Jos $K > 0$, →L13.
- L22. [Sarjatyypin valinta.] (Jos halutaan (R,r)-sarjan kertoimet, algoritmi loppuu tähän.)
COEFF(K)←COEFF(K)*VALUE(K), $K = 1, 2, \dots, N$.
(Jos halutaan (R,e)-sarjan kertoimet, algoritmi päättyy tähän.) COEFF(K)←COEFF(K)/VALUE(N),
 $K = 1, 2, \dots, N$. Algoritmi L päättyy tähän
(COEFF-kentissä on (E,e)-sarjan kertoimet). ■

Algoritmia L vastaava FORTRAN IV - aliohjelmaryhmä

Ohjelmoitaessa algoritmia L IBM 7094:lle FORTRAN IV-kielellä on algoritmiin vielä tehty lisäys: mikäli laskutoimitus on varmasti tarkka, se huomioidaan ohjelmassa vaihtamalla TYPE- kentän etumerkki miinukseksi. Tällöin tulevat kysymykseen tapaukset, joissa

- operaatio on negatointi (ei merkitystä, koska kerroin nollautuu joka tapauksessa)
- yhteen- tai vähennyslaskun operandi on nolla
- kertolaskun operandi tai jakaja on itseisarvoltaan yksi
- lopputuloksen arvo on nolla.

Tämä tieto saattaa olla tarpeellinen määrättäessä kumulatiivisen pyöristysvirheen jakautumaa, koska eräissä algoritmeissa (esimerkiksi matriisin kääntämisessä) näitä tapauksia on ratkaisevasti enemmän kuin satunnaisuuden perusteella voitaisiin olettaa.

Lisäksi ohjelman avulla voidaan 'määrätä' tietty laskutoimitus tarkaksi tai epätarkaksi, jolloin edellä luetelluilla tapauksilla ei ole vaikutusta TYPE-kentän etumerkkiin.

FORTRAN-ohjelma on ryhmä aliohjelmaa, jotka on liitetty yhdeksi monihaaraiseksi function-aliohjelmaksi, jotta tietuepinon tilanvarausta ei tarvitsisi määritellä uudelleen jokaisen algoritmin A alkeistoimituksen aikana. Seuraavassa esitellään tähän aliohjelmaryhmään, 'ryhmään L' kuuluvat aliohjelmat.

Nimikkeet QI, QJ, QK, ja QN ovat kokonaismuuttujia. Eräissä kutsuissa esiintyvällä kokonaismuuttujalla \mathbb{N} ei ole merkitystä; sen arvoksi tulee nolla.

Ryhmän L aliohjelmat

LBEGIN Kutsu: $\mathbb{N} = \text{LBEGIN}(\text{VALUE}, \text{TYPE}, \text{OPER1}, \text{OPER2}, \text{COEFF}, \text{M})$

Kutsun on esiinnettävä pääohjelmassa ennen muiden ryhmän L aliohjelmien kutsuja. Aliohjelmaryhmä saa tiedon taulukoiden VALUE(M), TYPE(M), ..., COEFF(M) sijainnista. Pinoindeksille I annetaan alkuarvo 1.

TYPE, OPER1 ja OPER2 ovat pääohjelmassa määriteltäviä kokonais-, VALUE ja COEFF reaalilukutaulukoita. Kokonaisluku M ilmoittaa kaikkien taulukoiden ulottuvuuden.

LNAME Kutsu: QI = LNAME(VAL)

Kaikki algoritmin A alkuarvot on ilmoitettava tällä kutsulla ryhmälle L. Alkuarvon, reaaliluvun VAL, nimikkeeksi asetetaan QI. Alkuarvo katsotaan epätarkaksi, so. siihen liittyy pyöristysvirhe (TYPE(QI):ksi tulee +1).

LNAMEX Kutsu: QI = LNAMEX(VAL)

Kuten LNAME, mutta alkuarvo katsotaan tarkaksi (TYPE(QI):ksi tulee -1).

LNEG Kutsu: QI = LNEG(QJ)

Muuttuja, jonka nimike on QJ, negatoidaan. Negaation nimikkeeksi asetetaan QI.

LADD Kutsu: $QI = LADD(QJ, QK)$
Muuttujat, joiden nimikkeet ovat QJ ja QK, lasketaan yhteen. Tuloksen nimike on QI.

LADDX Kutsu: $QI = LADDX(QJ, QK)$
Kuten LADD, mutta laskutoimitus katsotaan aina tarkaksi.

LADDN Kutsu: $QI = LADDN(QJ, QK)$
Kuten LADD, mutta laskutoimitus katsotaan aina epätarkaksi.

LSUB Kutsu: $QI = LSUB(QJ, QK)$
Muuttuja, jonka nimike on QK, vähennetään muuttujasta, jonka nimike on QJ. Erotuksen nimike on QI.

LSUBX Kutsu: $QI = LSUBX(QJ, QK)$
Kuten LSUB, mutta laskutoimitus katsotaan aina tarkaksi.

LSUBN Kutsu: $QI = LSUBN(QJ, QK)$
Kuten LSUB, mutta laskutoimitus katsotaan aina epätarkaksi.

LMUL Kutsu: $QI = LMUL(QJ, QK)$
Muuttujat, joiden nimikkeet ovat QJ ja QK, kerrotaan keskenään. Tulon nimike on QI.

LMULX Kutsu: $QI = LMULX(QJ, QK)$
Kuten LMUL, mutta laskutoimitus katsotaan aina tarkaksi.

LMULN Kutsu: $QI = LMULN(QJ, QK)$
Kuten LMUL, mutta laskutoimitus katsotaan aina epätarkaksi.

LDIV Kutsu: $QI = LDIV(QJ, QK)$
Suoritetaan jakolasku. QJ on jaettavan, QK jakajan ja QI osamäärän nimike.

LDIVX Kutsu: $QI = LDIVX(QJ, QK)$
Kuten LDIV, mutta laskutoimitus katsotaan aina tarkaksi.

LDIVN Kutsu: $QI = LDIVN(QJ, QK)$
Kuten LDIV, mutta laskutoimitus katsotaan aina epätarkaksi.

LEND Kutsu: $K = LEND(QN)$
Muuttujan, jonka nimike on QN, (R,r)-sarjan kertoimet lasketaan kenttiin COEFF(1), COEFF(2), ..., COEFF(QN). $COEFF(QN+1) = \dots = COEFF(M) = 0$.

LREL Kutsu: $K = LREL(COEFIC, M)$
Edeltävässä LEND-käskyssä mainitun muuttujan (nimike QN) (E,e)-sarjan kertoimet lasketaan kenttiin COEFIC(1), COEFIC(2), ..., COEFIC(QN).

COEFIC on pääohjelmassa määritelty reaalityyppinen taulukko, jonka ulottuvuus on M. Voi olla myös $COEFIC = COEFF$.

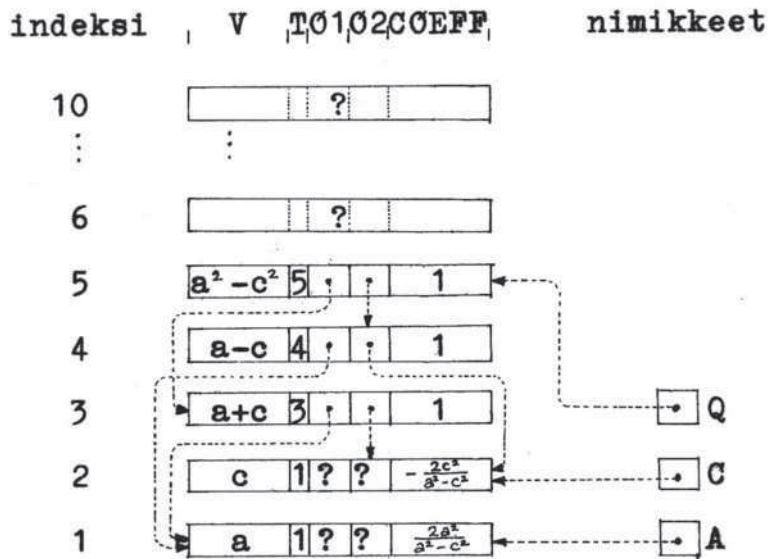
LABS Kutsu: $K = LABS(COEFIC, M)$
Kuten LREL, mutta COEFIC-taulukkoon tulevat (R,e)-sarjan kertoimet.

Jos LEND-kutsua seuraa sekä LREL- että LABS-kutsu, ei näistä ensimmäiselle saa olla $COEFIC = COEFF$.

Itse aliohjelmaryhmä L on esitetty liitteessä. Esimerkkinä sen käytöstä ohjelmoimme pienalgoritmin $a^2 - c^2 = (a+c) \cdot (a-c)$. Olkoon muuttujan VA arvo a ja muuttujan VC arvo c. Tällöin saamme $a^2 - c^2$:n (E,e)-sarjan kertoimet kenttiin COEFF(1), COEFF(2), ..., COEFF(Q) esimerkiksi ohjelmalla

```
INTEGER T(10), O1(10), O2(10), A, C, Q
REAL V(10), COEFF(10)
I = LBEGIN(V, T, O1, O2, COEFF, 10)
A = LNAME(VA)
C = LNAME(VC)
Q = LMUL(LADD(A, C), LSUB(A, C))
I = LEND(Q) + LREL(COEFF, 10) .
```

Aliohjelmaryhmän L muodostama rakenne arvoineen ohjelman suorituksen jälkeen on esitetty kuvassa 6, kysymysmerkillä merkityjä kenttiä ohjelma ei ole käsitellyt.



kuva 6

Yksikköhäiriön menetelmä

Algoritmin L varjopuolena on, että sen käyttö vaatii paljon muistitilaa, koska jokaiselle algoritmin välitulokselle on muodostettava oma tietuensa. Tilan säästämiseksi voidaan käyttää professori Tienarin esittämää 'yksikköhäiriön menetelmää', joka puolestaan käyttää runsaasti enemmän koneaikaa, koska algoritmi on tällöin laskettava läpi kerran mutakin yksittäistä pyöristysvirhettä kohden. Yksikköhäiriön menetelmää varten ei voida luoda ryhmää L vastaavia yleisiä aliohjelmia.

Yksikköhäiriön menetelmän periaatteena on, että algoritmin tuloksen q_N 'tarkka' arvo lasketaan kaksotarkkaa aritmetiikkaa käyttäen, jolloin voidaan asettaa yhtälössä (5.2) $E_N \approx 0$. Tämän jälkeen annetaan vuorollaan kullekin yksittäiselle suhteelliselle pyöristysvirheelle suuruusluokkaa b^{-t} oleva arvo. Kun algoritmi lasketaan 'vuorossa' oleva vir-

he e_i huomioon ottaen uudelleen läpi, saadaan tulokseksi q_N^* . Koska muut yksittäiset virheet ovat nollia, saamme yhtälön (5.2) muotoon

$$E_N = a_{N,i} e_i \quad (7.2)$$

Tiedämme, että $E_N = (q_N^* - q_N) / q_N$. Samoin tiedämme e_i :n suuruuden, joten kaavan (7.2) perusteella voimme ratkaista $a_{N,i}$:n kaavasta

$$a_{N,i} = \frac{q_N^* - q_N}{q_N e_i} \quad (7.3)$$

Vastaavalla tavalla saadaan $a_{N,i}$ ratkaistuksi kaikilla i :n arvoilla.

Esimerkkinä yksikköhäiriön menetelmästä ohjelmoimme jälleen pienalgoritmin $a^2 - c^2 = (a+c) \cdot (a-c)$. ER-muuttujan arvona on suuruusluokkaa b^{-t} oleva luku, ja virhekertoimet tulevat muuttujien COEFF(1), ..., COEFF(5) arvoiksi (haluttaessa ne voitaisiin tulostaa heti kun kukin niistä on laskettu). Taulukko E ajatellaan valmiiksi nollatuksi.

```
REAL COEFF(5)
DOUBLE PRECISION A,C,QA,QC,QX,QN,E(5)
QX = (A+B)*(A-C)
QXER = QX*ER
DO 10 I = 1,5
E(I) = ER
QA = A*(1.+E(1))
QC = C*(1.+E(2))
QN = (QA+QC)*(1.+E(3))*(QA-QC)*(1.+E(4))*(1.+E(5))
COEFF(I) = (QN-QX)/QXER
10 E(I) = 0.
```

Algoritmin L ja yksikköhäiriön menetelmän antamat tulokset ovat likimain yhtä tarkkoja tarkkuuden vähentyessä laskutoimitusten lukumäärän kasvassa. Mikäli tämä tarkkuus ei riitä, voidaan algoritmi L ohjelmoida käyttämään kaksioistarkkaa aritmetiikkaa, jolloin lasketut arvot ovat ratkaisevasti tarkempia.

8. PIENALGORITMIT $a^2 - c^2$:N LASKEMISEKSI

Esimerkkinä pyöristysvirheen Taylor-sarjan käytöstä pyrimme selvittämään sen avulla, kumpi lausekkeen $a^2 - c^2$ kahdesta mahdollisesta laskutavasta,

$$1. \quad a^2 - c^2 = (a+c) \cdot (a-c) \quad (8.1)$$

vai

$$2. \quad a^2 - c^2 = a \cdot a - c \cdot c \quad (8.2)$$

on edullisempi.

Algoritmien yksinkertaisuuden vuoksi voimme laskea niiden Taylor-sarjat luvussa 5 esitetyllä tavalla. Tällöin

$$\begin{aligned} 1. & \quad ((a^* + c^*)^* \cdot (a^* - c^*)^*)^* \\ & = \{[a(1+e_a) + c(1+e_c)](1+e_1) [a(1+e_a) - c(1+e_c)](1+e_2)\}(1+e_3) \\ & = a^2 - c^2 + 2a^2 e_a - 2c^2 e_c + (a^2 - c^2) e_1 + (a^2 - c^2) e_2 + (a^2 - c^2) e_3 + \langle e^2 \rangle \quad (8.3) \\ & = (a^2 - c^2) \left(1 + \frac{2a^2}{a^2 - c^2} e_a - \frac{2c^2}{a^2 - c^2} e_c + e_{s_1} + e_{s_2} + e_{s_3} + \langle e^2 \rangle\right) \\ & = (a - c)(1 + E_1) \end{aligned}$$

ja

$$\begin{aligned} 2. & \quad ((a^* \cdot a^*)^* - (c^* \cdot c^*)^*)^* \\ & = \{[a(1+e_a) a(1+e_a)](1+e_4) - [c(1+e_c) c(1+e_c)](1+e_5)\}(1+e_6) \\ & = a^2 - c^2 + 2a^2 e_a - 2c^2 e_c + a^2 e_4 - c^2 e_5 + (a^2 - c^2) e_6 + \langle e^2 \rangle \quad (8.4) \\ & = (a^2 - c^2) \left(1 + \frac{2a^2}{a^2 - c^2} e_a - \frac{2c^2}{a^2 - c^2} e_c + \frac{a^2}{a^2 - c^2} e_4 - \frac{c^2}{a^2 - c^2} e_5 + e_{s_6} + \langle e^2 \rangle\right) \\ & = (a^2 - c^2)(1 + E_2) \end{aligned}$$

missä e_a ja e_c ovat alkuarvojen a ja c pyöristysvirheitä, e_3, e_4 ja e_5 tulon sekä e_1, e_2 ja e_6 summan tai erotuksen pyöristysvirheitä. Edellä (kuva 6) totesimme tapauksessa 1 algoritmin L johtavan samaan Taylor-kehitelämään kuin (8.3).

Oletamme, että $|a| > |c|$, mikä ei ole oleellinen rajoitus. Jos tarkastelemme suurimpia mahdollisia virheitä, saamme

$$\max(E_1) \approx \left(\frac{2a^2}{|a^2-c^2|} + \frac{2c^2}{|a^2-c^2|} + 3 \right) \cdot \max(e), \quad (8.5)$$

$$\max(E_2) \approx \left(\frac{3a^2}{|a^2-c^2|} + \frac{3c^2}{|a^2-c^2|} + 1 \right) \cdot \max(e), \quad (8.6)$$

missä $\max(e)$ tarkoittaa suurinta mahdollista yksittäistä pyöristysvirhettä. Esimerkiksi pyöristäväsä aritmetiikassa, jossa pyöristykset suoritetaan laskutoimituksen jälkeen, on yhtälöiden (3.17) ja (4.3) perusteella $\max(e) = \frac{1}{2}ub$.

Yhtälöiden (8.5) ja (8.6) perusteella $\max(E_1) < \max(E_2)$ eli tapa 1 on edullisempi, kun

$$2|a^2 - c^2| < a^2 + c^2 \quad (8.7)$$

eli, koska $|a| > |c|$, kun

$$\left| \frac{a}{c} \right| < \sqrt{3} \approx 1.7302. \quad (8.8)$$

Jos pyrimme selvittämään, kumpi laskutavoista on suuremmalla todennäköisyydellä edullinen, meidän on tarkasteltava virheiden odotusarvoja ja variansseja. Kaavojen (5.8) ja (5.9) perusteella

$$E(E_1) \approx 2\mu_A + 2\mu_S + \mu_T, \quad (8.9)$$

$$E(E_2) \approx 2\mu_A + \mu_S + \mu_T, \quad (8.10)$$

$$D^2(E_1) \approx 4 \frac{a^4 + c^4}{(a^2 - c^2)^2} \sigma_A^2 + 2\sigma_S^2 + \sigma_T^2, \quad (8.11)$$

$$D^2(E_2) \approx 4 \frac{a^4 + c^4}{(a^2 - c^2)^2} \sigma_A^2 + \sigma_S^2 + \frac{a^4 + c^4}{(a^2 - c^2)^2} \sigma_T^2. \quad (8.12)$$

Yksinkertaistamme tehtävää olettamalla, että aritmetiikka on pyöristävä, jolloin $E(E_1) = E(E_2) = 0$. Tapa 1 on tällöin edullisempi, kun $D^2(E_1) < D^2(E_2)$ eli

$$(a^2 - c^2)^2 \sigma_S^2 < 2a^2 c^2 \sigma_T^2. \quad (8.13)$$

Kun $|a| \approx |c|$, on $a^2 - c^2 \approx 0$ ja tapa 1 on edullisem-

pi. Muulloin, ellei $|a| \gg |c|$, on $\sigma_s^2 \approx \sigma_r^2$. Tällöin tapa 1 on kaavan (8.13) mukaan edullisempi, mikäli

$$\left| \frac{a}{c} \right| < \sqrt{2+\sqrt{3}} \approx 1.932 \quad . \quad (8.14)$$

Kun a :n ja c :n eksponenttien erotus on t , on

$$(a^* \pm c^*)^* = a^* \quad . \quad (8.15)$$

Myös eksponenttien erotuksen arvolla t yhtälö (8.15) on voimassa, kunhan c :n mantissan itseisarvo $< \frac{1}{2}$. Koska a :n mantissan itseisarvo on kaavan (2.6) mukaan < 1 , voimme todeta yhtälön (8.15) olevan voimassa aina, jos

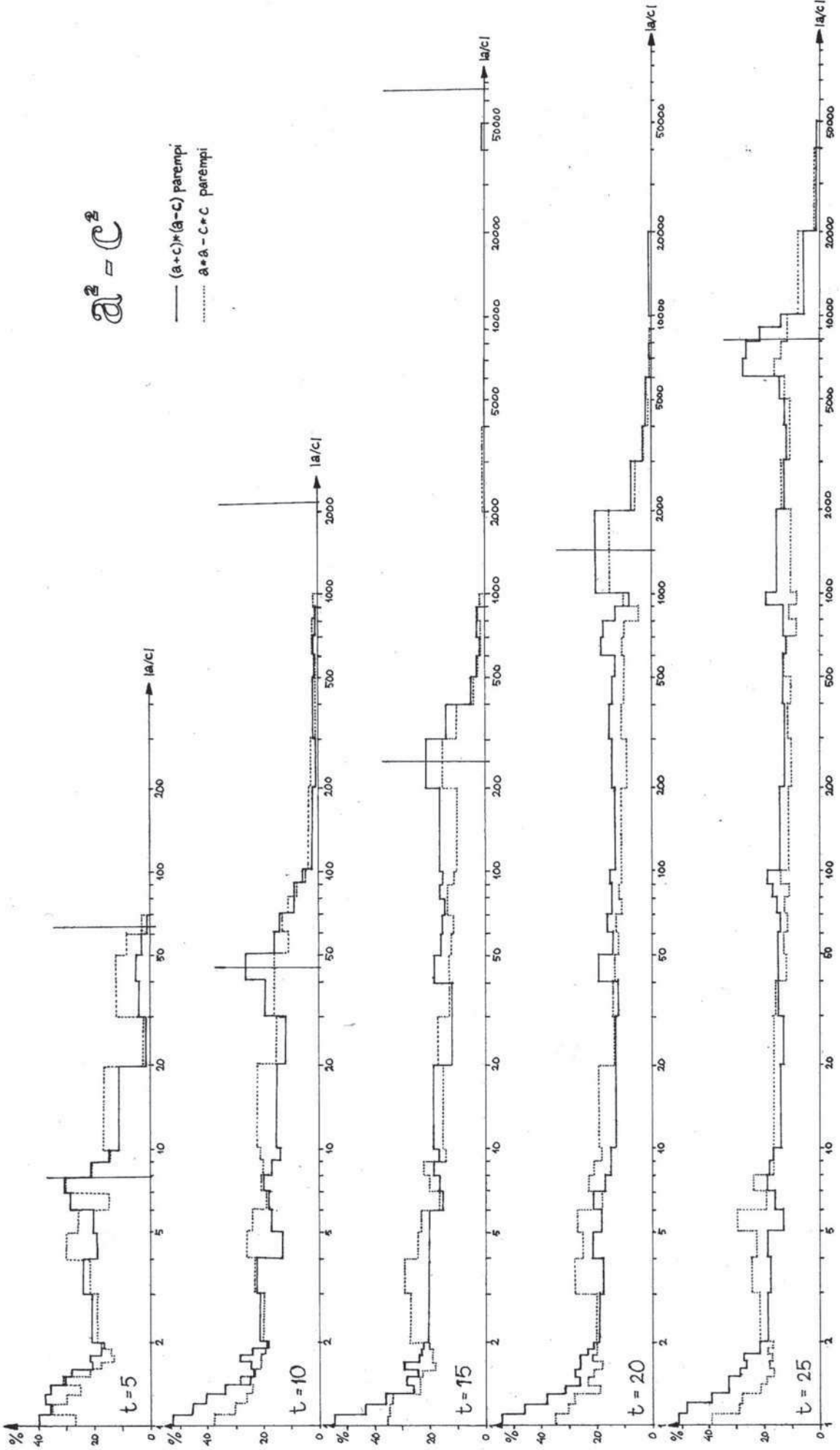
$$\left| \frac{a}{c} \right| \approx \left| \frac{a^*}{c^*} \right| > \frac{1}{\frac{1}{2} \cdot b^{-t}} = 2b^t \quad . \quad (8.16)$$

Laskutapa 1 supistuu siis tällöin muotoon $a \cdot a$. Tavassa 2 voimme vastaavasti korvata c :n nollalla, jos

$$\left| \frac{a}{c} \right| \approx \sqrt{\frac{(a^* \cdot a^*)^*}{(c^* \cdot c^*)^*}} > \sqrt{2b^t} = \sqrt{2} \cdot b^{t/2} \quad . \quad (8.17)$$

Koska $2 \cdot b^t > \sqrt{2} \cdot b^{t/2}$, kun $b \geq 2$ ja $t \geq 1$, riittää ehto (8.16) molemmissa tavoissa takaamaan, että c :llä ole vaikutusta lopputulokseen. Kaavoissa (8.3) ja (8.4) tämä merkitsee, että c, e_1, e_2 ja e_6 voidaan korvata nollalla, jolloin laskutavat päätyvät samaan tulokseen eikä siis kumpikaan ole toistaan edullisempi.

Saatujen tulosten testaamiseksi laskettiin $a^2 - c^2$:n arvoja satunnaisilla a :n ja c :n arvoilla. Tulosten luetteloimiseksi ne jaettiin $|a/c|$:stä riippuviin luokkiin, välillä $[1,2)$ toiseksi merkitsevimmän ja välillä $[2,100000)$ merkitsevimmän (desimaali)numeron perusteella. Kussakin luokassa suoritettiin noin 300 laskutoimitusta kummallakin tavalla. Kaik-



kuva 7

ki laskutoimitukset suoritettiin binääriaritmetiikassa ($b = 2$) viidellä erilaisella numeroiden lukumäärällä ($t = 5, 10, 15, 20$ ja 25).

Testin tulos on histogrammana kuvassa 7, jossa abskissana on $|a/c|$:n arvo ja oordinaattana prosenttinen osuus suoritetuista laskutoimituksista. Eheä viiva ilmaisee, monessako prosentissa tapa 1 oli parempi ja katkoviiva vastaavan prosenttimäärän tavasta 2. Lopuissa tapauksissa tulos oli molemmilla tavoilla laskettuna sama. $|a/c|$:n edellä mainitut arvot $\sqrt{2} \cdot b^{t/2}$ ja $2b^t$ on merkitty histogrammeihin pystyviivoilla.

Voimme todeta tavan 1 todella olevan edullisempi, kun $|a/c| \leq 2$. Mitä lähempänä $|a|$ ja $|c|$ ovat toisiaan, sitä todennäköisempää on, että laskutavat 1 ja 2 johtavat eri tukoksiin. Kun $|a/c| > 2$, ei suoritettun testin perusteella voida sanoa, että tapa 2 olisi merkittävästi parempi.

Samaan tulokseen päätyminen todennäköisyys pysyy $|a/c|$:n kasvaessa likimain vakiona aina rajaan $\sqrt{2} \cdot b^{t/2}$ asti. Tämän jälkeen alkaa olla yhä todennäköisempää, että laskutavat antavat saman tuloksen. Kuten kaavan (8.16) perusteella voitiin odottaa, ei eriäviä tuloksia saatu, kun $|a/c|$ oli $> 2b^t$.

9. HORNER-SHEMA

Taylor-kehityksen analyttinen määrittäminen

Horner-shema on algoritmi, jolla lasketaan polynomin

$$p = a_0 x^N + a_1 x^{N-1} + \dots + a_{N-1} x + a_N, \quad a_0 \neq 0, \quad (9.1)$$

arvo pisteessä x . Se kirjoitetaan tavallisesti muotoon [4]

$$\begin{cases} q_0 = a_0 \\ q_n = a_n + xq_{n-1}, \quad n = 1, \dots, N. \end{cases} \quad (9.2)$$

Kun hajotamme algoritmin (9.2) yksittäisiin laskutoimituksiin, saamme sen muotoon

$$\begin{cases} q_{0,2} = a_0 \\ q_{n,1} = xq_{n-1,2} \\ q_{n,2} = a_n + q_{n,1}, \quad n = 1, \dots, N. \end{cases} \quad (9.3)$$

Sovellamme tähän algoritmiin analyttistä menetelmää kumulatiivisen pyöristysvirheen Taylor-sarjan kertoimien laskemiseksi.

Kaavan (6.12) perusteella saamme

$$\begin{cases} R_{n,1}^{(i)} = q_{n-1,2} x e_x + x R_{n-1,2}^{(i)} + q_{n,1} e_{n,1} \\ R_{n,2}^{(i)} = a_n e_n + R_{n,1}^{(i)} + q_{n,2} e_{n,2}, \end{cases} \quad (9.4)$$

missä e_x on x :n ja e_n a_n :n pyöristysvirhe sekä $e_{n,1}$ ja $e_{n,2}$ n :nnen iteraatiokierron kerto- ja yhteenlaskun pyöristysvirheet.

Sijoittamalla yhtälöistä (9.4) edellinen jälkimmäiseen saamme ensimmäisen asteen differenssiyhtälön

$$R_{n,2}^{(4)} = xR_{n-1,2}^{(4)} + xq_{n-1,2}e_x + a_n e_n + xq_{n-1,1}e_{n,1} + q_{n,2}e_{n,2}, \quad (9.5)$$

jonka alkuarvona on $R_{0,2}^{(4)} = a_0 e_0$ eli a_0 :n absoluuttinen pyöristysvirhe.

Muotoa

$$x_n = ax_{n-1} + b_n, \quad x_0 = b_0 \quad (9.6)$$

olevan differenssiyhtälön ratkaisu on [4]

$$x_n = \sum_{i=0}^n a^{n-i} b_i, \quad (9.7)$$

joten saamme differenssiyhtälön (9.5) ratkaisuksi

$$R_{n,2}^{(4)} = \sum_{i=1}^n x^{n-i+1} q_{i-1,2} e_x + \sum_{i=0}^n x^{n-i} a_i e_i + \sum_{i=1}^n x^{n-i+1} q_{i-1,1} e_{i,1} + \sum_{i=1}^n x^{n-i} q_{i,2} e_{i,2}, \quad n = 2, \dots, N. \quad (9.8)$$

Vastaavasti saadaan

$$\begin{cases} R_{1,1}^{(4)} = xq_{0,2}e_x + xa_0e_0 + xq_{0,1}e_{1,1} \\ R_{n,1}^{(4)} = \sum_{i=1}^n x^{n-i+1} q_{i-1,1} e_x + \sum_{i=1}^{n-1} x^{n-i} a_i e_i + \sum_{i=1}^n x^{n-i+1} q_{i-1,2} e_{i,2} \\ \quad + \sum_{i=1}^{n-1} x^{n-i} q_{i,1} e_{i,1}, \quad n = 2, \dots, N. \end{cases} \quad (9.9)$$

Tarkastelemme erityisesti algoritmin lopputuloksen $p = q_{N,2}$ pyöristysvirhettä $R_{N,2}$. Otamme käyttöön merkinnät

$$\begin{cases} \alpha_i = x^{N-i} a_i \\ \beta_i = \sum_{j=0}^i \alpha_j \\ \gamma_i = \sum_{j=0}^i (N-j) \alpha_j \end{cases} \quad (9.10)$$

Merkinnällä α_i tarkoitamme siis polynomin p $(N-i)$:nnen asteen termiä ja β_i :llä $(N-i)$:nnen ja sitä korkeamman asteen termien summaa. Kun toteamme, että

$$q_{i,2} = \sum_{j=0}^i a_j x^{i-j}, \quad (9.11)$$

saamme β_i :lle myös lausekkeen

$$\beta_i = x^{N-i} q_{i,2}. \quad (9.12)$$

Erityisesti

$$\beta_N = p \quad (9.13)$$

Polynomien p derivaatan p' arvo on

$$p' = Na_0 x^{N-1} + (N-1)a_1 x^{N-2} + \dots + 2a_{N-2}x + a_{N-1}. \quad (9.14)$$

Huomaamme, että γ_i on polynomien $x p'$ $(N-i)$:nnen ja sitä korkeamman asteen termien summa, erityisesti

$$\gamma_N = x p' \quad (9.15)$$

Voimme lausua γ_i :n myös muodossa

$$\gamma_i = \sum_{j=0}^i (N-i)\alpha_j + \sum_{j=0}^i (i-j)\alpha_j = (N-i)\beta_j + \sum_{j=0}^{i-1} \beta_j, \quad (9.16)$$

erityisesti

$$\gamma_N = \gamma_{N-1} = \sum_{i=0}^{N-1} \beta_i = \sum_{i=1}^{N-1} x^{N-i+1} q_{i-1,2}. \quad (9.17)$$

Saamme yhtälön (9.8) yhtälöiden (9.10), (9.12) ja (9.17) avulla muotoon

$$R_{N,2}^{(1)} = \gamma_N e_x + \sum_{i=0}^N \alpha_i e_i + \sum_{i=1}^N \beta_{i-1} e_{i,1} + \sum_{i=1}^N \beta_i e_{i,2}, \quad N \geq 2, \quad (9.18)$$

josta yhtälön (6.4) perusteella saamme polynomien p (E, e) -sarjaksi

$$E_{N,2}^{(1)} = \frac{\gamma_N}{p} e_x + \sum_{i=0}^N \frac{\alpha_i}{p} e_i + \sum_{i=1}^N \frac{\beta_{i-1}}{p} e_{i,1} + \sum_{i=1}^N \frac{\beta_i}{p} e_{i,2}, \quad N \geq 2. \quad (9.19)$$

Kaavojen (5.8) ja (5.9) mukaan saamme $E_{N,2}^{(1)}$:n odotusarvon ja varianssin muotoon

$$E(E_{N,2}^{(1)}) = \frac{1}{p} (\gamma_N + \sum_{i=0}^N \alpha_i) \mu_A + \sum_{i=1}^N \frac{\beta_i}{p} \mu_S + \sum_{i=1}^N \frac{\beta_{i-1}}{p} \mu_T, \quad (9.20)$$

$$D^2(E_{N,2}^{(1)}) = \frac{1}{p^2} (\gamma_N^2 + \sum_{i=0}^N \alpha_i^2) \sigma_A^2 + \sum_{i=1}^N \frac{\beta_i^2}{p^2} \sigma_S^2 + \sum_{i=1}^N \frac{\beta_{i-1}^2}{p^2} \sigma_T^2. \quad (9.21)$$

Kun oletamme alkuarvot tarkoiksi ja laskutoimitukset suoritettavaksi kiinteän pilkun aritmetiikkaa käyttäen, jolloin yhteenlaskut ovat tarkkoja,

on yhtälössä (9.18) $e_x = e_i = e_{i,2} = 0$, $i = 0, \dots, N$.
Kaavojen (6.17) ja (9.12) perusteella saamme yhtälön (9.18) tällöin muotoon

$$R_{N,2}^{(1)} = \sum_{i=1}^N \beta_{i-1} \frac{r_{i,1}}{q_{i,1}} = \sum_{i=1}^N x^{N-i+1} q_{i-1,2} \frac{r_{i,1}}{x q_{i-1,2}} = \sum_{i=1}^N x^{N-i} r_{i,1}, \quad (9.22)$$

mikä on sama kuin Henricin saama tulos. [4].

Toisen asteen kertoimet

Myös toisen asteen kertoimet Horner-sheman kumulatiivisen pyöristysvirheen Taylor-kehityksessä voidaan johtaa analyttisesti. Kaavojen (6.13) ja (9.3) perusteella

$$\begin{cases} R_{n,1}^{(2)} = q_{n-1,2} \cdot 0 + x R_{n-1,2}^{(2)} + x e_x R_{n-1,2}^{(1)} + q_{n-1,2} x e_x e_{n,1} + x R_{n-1,2}^{(1)} e_{n,1} \\ R_{n,2}^{(2)} = 1 \cdot 0 + 1 \cdot R_{n,1}^{(2)} + 0 + a_n e_n e_{n,2} + R_{n,1}^{(1)} e_{n,2} \end{cases} \quad (9.21)$$

Sijoittamalla ensimmäisen yhtälön toiseen saamme ensimmäisen asteen differenssiyhtälön

$$R_{n,2}^{(2)} = x R_{n-1,2}^{(2)} + x e_x R_{n-1,2}^{(1)} + x q_{n-1,2} e_x e_{n,1} + x R_{n-1,2}^{(1)} e_{n,1} + a_n e_n e_{n,2} + R_{n,1}^{(1)} e_{n,2}. \quad (9.22)$$

Alkuarvon $R_{0,2}^{(2)} = 0$ avulla saamme sen ratkaisuksi

$$R_{n,2}^{(2)} = \sum_{i=1}^n x^{n-i+1} e_x R_{i-1,2}^{(1)} + \sum_{i=1}^n x^{n-i+1} q_{i-1,2} e_x e_{i,1} + \sum_{i=1}^n x^{n-i+1} R_{i-1,2}^{(1)} e_{i,1} + \sum_{i=1}^n x^{n-i} a_i e_i e_{i,2} + \sum_{i=1}^n x^{n-i} R_{i,1}^{(1)} e_{i,2}, \quad n = 1, \dots, N. \quad (9.23)$$

Sijoittamalla tähän yhtälöön kaavojen (9.8) ja (9.9) tulokset saamme

$$\begin{aligned} R_{n,2}^{(2)} = & a_0 x^n e_0 (e_x + e_{1,1}) + \sum_{i=2}^n x^{n-i+1} \left(\sum_{j=1}^{i-1} x^{i-j} q_{j-1,2} e_x + \sum_{j=0}^{i-1} x^{i-j-1} a_j e_j \right. \\ & \left. + \sum_{j=1}^{i-1} x^{i-j} q_{j-1,2} e_{j,1} + \sum_{j=1}^{i-1} x^{i-j-1} q_{j,2} e_{j,2} \right) (e_x + e_{i,1}) + \sum_{i=1}^n x^{n-i+1} q_{i-1,2} e_x e_{i,1} \\ & + \sum_{i=1}^n x^{n-i} a_i e_i e_{i,2} + \sum_{i=1}^n x^{n-i} \left(\sum_{j=1}^i x^{i-j+1} q_{j-1,2} e_x + \sum_{j=0}^{i-1} x^{i-j} a_j e_j \right. \\ & \left. + \sum_{j=1}^i x^{i-j+1} q_{j-1,2} e_{j,1} \right) e_{i,2} + \sum_{i=2}^n \sum_{j=1}^{i-1} x^{n-i} x^{i-j} q_{j,2} e_{j,2} e_{i,2}, \\ & n = 2, \dots, N. \end{aligned} \quad (9.24)$$

Erityisesti kaavojen (9.10) ja (9.12) perusteella

$$\begin{aligned}
 R_{N,2}^{(2)} = & \alpha_0 \mathbf{e}_0 (\mathbf{e}_x + \mathbf{e}_{i,1}) + \sum_{i=2}^N \left(\sum_{j=1}^{i-1} \beta_{j-1} \mathbf{e}_x + \sum_{j=0}^{i-1} \alpha_j \mathbf{e}_j + \sum_{j=1}^{i-1} \beta_{j-1} \mathbf{e}_{j,1} + \sum_{j=1}^{i-1} \beta_j \mathbf{e}_{j,2} \right) (\mathbf{e}_x + \mathbf{e}_{i,1}) \\
 & + \sum_{i=1}^N \beta_{i-1} \mathbf{e}_x \mathbf{e}_{i,1} + \sum_{i=1}^N \alpha_i \mathbf{e}_i \mathbf{e}_{i,2} + \sum_{i=1}^N \left(\sum_{j=1}^i \beta_{j-1} \mathbf{e}_x + \sum_{j=0}^{i-1} \alpha_j \mathbf{e}_j \right. \\
 & \left. + \sum_{j=1}^i \beta_{j-1} \mathbf{e}_{j,1} \right) \mathbf{e}_{i,2} + \sum_{i=2}^N \sum_{j=1}^{i-1} \beta_j \mathbf{e}_{j,2} \mathbf{e}_{i,2} \quad .
 \end{aligned} \quad (9.25)$$

Kaavan

$$\sum_{i=k}^N \sum_{j=k-1}^{i-1} c_j = \sum_{j=k-1}^{N-1} (N-j) c_j \quad (9.26)$$

avulla saamme yhtälön (9.25) muotoon

$$\begin{aligned}
 R_{N,2}^{(2)} = & \sum_{j=1}^{N-1} (N-j) \beta_{j-1} \mathbf{e}_x \mathbf{e}_x + \sum_{j=0}^N (N-j) \alpha_j \mathbf{e}_x \mathbf{e}_j + \sum_{i=1}^N [(N-1) \beta_{i-1} \\
 & + \sum_{j=1}^i \beta_{j-1}] \mathbf{e}_x \mathbf{e}_{i,1} + \sum_{i=1}^N [(N-1) \beta_i + \sum_{j=1}^i \beta_{j-1}] \mathbf{e}_x \mathbf{e}_{i,2} \\
 & + \sum_{i=1}^N \sum_{j=0}^{i-1} \alpha_j \mathbf{e}_j \mathbf{e}_{i,1} + \sum_{i=1}^N \sum_{j=0}^{i-1} \alpha_j \mathbf{e}_j \mathbf{e}_{i,2} + \sum_{i=1}^N \alpha_i \mathbf{e}_i \mathbf{e}_{i,2} \quad (9.27) \\
 & + \sum_{i=2}^N \sum_{j=1}^{i-1} \beta_{j-1} \mathbf{e}_{i,1} \mathbf{e}_{j,1} + \sum_{i=1}^N \sum_{j=1}^{i-1} \beta_j \mathbf{e}_{i,2} \mathbf{e}_{j,2} + \sum_{i=1}^N \sum_{j=1}^i \beta_{j-1} \mathbf{e}_{j,1} \mathbf{e}_{i,2} \\
 & + \sum_{i=2}^N \sum_{j=1}^{i-1} \beta_j \mathbf{e}_{i,1} \mathbf{e}_{j,2} \quad , \quad N \geq 2 \quad .
 \end{aligned}$$

Kaavan (9.10) perusteella

$$\begin{aligned}
 \sum_{j=1}^{N-1} (N-j) \beta_{j-1} & = \sum_{j=1}^{N-1} \sum_{k=0}^{j-1} (N-j) \alpha_k = \sum_{k=0}^{N-2} \frac{1}{2} (N-k-1) (N-k) \alpha_k \\
 & = \frac{1}{2} x^2 p'' \quad , \quad (9.28)
 \end{aligned}$$

missä p'' on polynomin p toinen derivaatta pisteessä x . Kun vielä käytämme kaavoja (9.16) ja (6.4), saamme

$$\begin{aligned}
 E_{n,2}^{(2)} = & \frac{1}{2p} x^2 p'' \mathbf{e}_x \mathbf{e}_x + \sum_{j=0}^N (N-j) \frac{\alpha_j}{p} \mathbf{e}_x \mathbf{e}_j + \sum_{i=1}^N \frac{\beta_{i-1}}{p} \mathbf{e}_x \mathbf{e}_{i,1} + \sum_{i=1}^N \frac{\beta_i}{p} \mathbf{e}_x \mathbf{e}_{i,2} \\
 & + \sum_{i=1}^N \sum_{j=0}^{i-1} \frac{\alpha_j}{p} \mathbf{e}_j \mathbf{e}_{i,1} + \sum_{i=1}^N \sum_{j=0}^{i-1} \frac{\alpha_j}{p} \mathbf{e}_j \mathbf{e}_{i,2} + \sum_{i=2}^N \sum_{j=1}^{i-1} \frac{\beta_{j-1}}{p} \mathbf{e}_{i,1} \mathbf{e}_{j,1} \quad (9.29) \\
 & + \sum_{i=2}^N \sum_{j=1}^{i-1} \frac{\beta_j}{p} \mathbf{e}_{i,1} \mathbf{e}_{j,2} + \sum_{i=1}^N \sum_{j=1}^{i-1} \frac{\beta_j}{p} \mathbf{e}_{j,1} \mathbf{e}_{i,2} + \sum_{i=2}^N \sum_{j=1}^{i-1} \frac{\beta_j}{p} \mathbf{e}_{i,2} \mathbf{e}_{j,2} \quad , \quad N \geq 2 \quad .
 \end{aligned}$$

Kaavojen (9.19) ja (9.29) antamat tulokset p:n (E,e)-sarjan ensimmäisen ja toisen asteen kertoimille on esitetty alla olevassa taulukossa, kun polynomin asteluku $N \geq 1$.

tekijä	kerroin*p	ehto tekijän indekseille	ao. termien lukumäärä
e_x	γ_N	-	1
e_i	α_i	$i \geq 0$	$N+1$
$e_{i,1}$	β_{i-1}	$i \geq 1$	N
$e_{i,2}$	β_i	$i \geq 1$	N
$e_x e_x$	$\frac{1}{2} x^2 p^n$	-	1 (0, jos $N=1$)
$e_x e_i$	$(N-1)\alpha_i$	$i \geq 0$	$N+1$
$e_x e_{i,1}$	γ_{i-1}	$i \geq 1$	N
$e_x e_{i,2}$	γ_i	$i \geq 1$	N
$e_i e_{j,1}$	α_i	$j > i \geq 0$	$\frac{1}{2}(N^2 + N)$
$e_{i,1} e_{j,2}$	α_i	$j > 1 \geq 0$ tai $i=j \geq 1$	$\frac{1}{2}(N^2 + 3N)$
$e_{i,1} e_{j,1}$	β_{i-1}	$j > i \geq 1$	$\frac{1}{2}(N^2 - N)$
$e_{i,1} e_{j,2}$	β_j	$i > j \geq 1$	$\frac{1}{2}(N^2 - N)$
$e_{i,1} e_{j,2}$	β_{i-1}	$j \geq i \geq 1$	$\frac{1}{2}(N^2 + N)$
$e_{i,2} e_{j,2}$	β_i	$j > i \geq 1$	$\frac{1}{2}(N^2 - N)$

Voimme todeta ensimmäisen asteen termejä olevan kaikkiaan $3N+2$ ja toisen asteen termejä $3N^2+4N+2$ kappaletta eli toisen asteen termien lukumäärä on suurilla N:n arvoilla noin N-kertainen verrattuna ensimmäisen asteen termien lukumäärään.

Ensimmäisen ja toisen asteen kertoimet edustavat samaa suuruusluokkaa, mutta e^2 on suuruusluokkaa $b^{-t}e$. Jotta toisen asteen termien vaikutus olisi samaa suuruusluokkaa kuin ensimmäisen asteen termien, olisi siis N:n oltava suuruusluokkaa b^t . Tällöin pyöristysvirheet ovat kuitenkin niin suuria, ettei polynomin arvon laskeminen t numeron tarkkuudella yleensä enää ole mielekästä tuloksen epätarkkuuden johdosta.

Esimerkkinä mainitusta epätarkkuudesta tarkastelemme polynomia $p = x^N$, missä $N = b^t$. Koska $a_0 = 1$ ja $a_i = 0$, $i = 1, \dots, N$, on $e_{1,1} = 0$ ja $e_i = e_{i,2} = 0$, $i = 1, \dots, N$. Kaavojen (9.20), (3.23) ja (4.9) perusteella on $E_{N,2}^{(1)}$:n odotusarvo katkaisevassa aritmetiikassa

$$E(E_{N,2}^{(1)}) = \frac{1}{p} \left(\gamma_N \mu_\lambda + \sum_{i=2}^N \beta_{i-1} \mu_\tau \right) = \frac{1}{p} (N x^N \mu_\lambda + (N-1) x^N \mu_\tau) \\ \approx 2N \frac{u(1-b)}{2 \ln b} = \frac{1-b}{\ln b} .$$

Toisistaan riippumattomien satunnaismuuttujien ξ ja η tulon odotusarvolle voimassa olevan yhtälön [2]

$$E(\xi\eta) = E(\xi)E(\eta) \quad (9.30)$$

avulla saamme $E_{N,2}^{(2)}$:n odotusarvoksi katkaisevassa aritmetiikassa

$$E(E_{N,2}^{(2)}) = \frac{1}{p} \left(\frac{1}{2} x^2 p'' \mu_\lambda^2 + \sum_{i=2}^N \gamma_{i-1} \mu_\lambda \mu_\tau + \sum_{i=3}^N \sum_{j=2}^{i-1} \beta_{j-1} \mu_\tau^2 \right) \\ = \frac{1}{p} \left(\frac{1}{2} x^2 N(N-1) x^{N-2} \mu_\lambda^2 + (N-1) N x^N \mu_\lambda \mu_\tau + \frac{(N-2)(N-1)}{2} x^N \mu_\tau^2 \right) \\ \approx 2N^2 \left[\frac{u(1-b)}{2 \ln b} \right]^2 = \frac{1}{2} \left[\frac{1-b}{\ln b} \right]^2 .$$

Alla olevaan taulukkoon on laskettu saatujen odotusarvojen arvoja kantaluvun vaihdellessa.

b	$E(E_{N,2}^{(1)})$	$E(E_{N,2}^{(2)})$
2	-1.44	1.04
8	-3.37	5.67
10	-3.91	7.64
16	-5.41	14.6
64	-15.1	115

Kuten odotimme, edustavat $E(E_{N,2}^{(1)})$ ja $E(E_{N,2}^{(2)})$ samaa suuruusluokkaa, ja niiden arvot ovat niin suuria, että virheen odotusarvo varsinkin suurilla kantaluvun arvoilla on moninkertainen itse polynomien arvoon verrattuna.

Polynomille $p = x^N$, missä $N = \frac{1}{2}b^{t-1}$, saamme vastaavaksi taulukoksi

b	$E(E_{N,2}^{(1)})$	$E(E_{N,2}^{(2)})$
2	-0.36	0.07
8	-0.21	0.022
10	-0.19	0.019
16	-0.17	0.014
64	-0.12	0.0045

Toisen asteen termien vaikutus on jo selvästi ensimmäisen asteen termien vaikutusta pienempi, mutta virheen odotusarvo on vielä varsin suuri.

Saatu tulos vahvistaa yhtälössä (5.4) esitettyä käsitystä, jonka mukaan toisen ja korkeamman asteen termeillä ei ole käytännössä merkitystä kumulatiivista pyöristysvirhettä tarkasteltaessa.

Nollakohtien lähekkäisyyden vaikutus kertoimiin

On tunnettua, että polynomien nollakohtien suhteellinen lähekkäisyys vaikuttaa heikentävästi polynomien arvon laskutarkkuuteen. Tutkimme seuraavassa tätä ilmiötä Taylor-kehityksen (9.19) kertoimien avulla.

Polynomien p (9.1) kerroin a_i voidaan lausua nollakohtien z_j , $j = 1, \dots, N$, avulla muodossa

$$a_i = a_0 (-1)^i \sum_{j=1}^K (z_{j_1} z_{j_2} \dots z_{j_i}), \quad i = 1, \dots, N, \quad (9.31)$$

missä summattavien tulojen tekijöinä ovat kaikki mahdolliset i nollakohdan kombinaatiot, jolloin $K = N!/(i!(N-i)!)$.

Tarkastelemme nyt polynomia \bar{p} , jonka nollakohdat ovat $\bar{z}_j = z_j + M$, $j = 1, \dots, N$, ja jolle

$$\bar{p}(\bar{x}) = \bar{p}(x+M) = p(x) \quad . \quad (9.32)$$

Tällöin

$$\begin{aligned} \bar{p}(\bar{x}) &= \bar{a}_0 [(x+M) - (z_1+M)] \cdots [(x+M) - (z_N+M)] \\ &= \bar{a}_0 (x-z_1) \cdots (x-z_N) = \frac{\bar{a}_0}{a_0} p(x), \end{aligned} \quad (9.33)$$

joten

$$\bar{a}_0 = a_0 \quad (9.34)$$

ja kaavan (9.31) mukaan

$$\bar{a}_i = a_0 (-1)^i \sum_{j=1}^k [(z_{j_1}+M) \cdots (z_{j_i}+M)] \quad (9.35)$$

Kun M kasvaa itseisarvoltaan riittävän isoksi, ovat \bar{z}_j :t, \bar{x} ja M samanmerkkisiä. $|M|$:n kasvaessa edelleen kasvaa myös

$$\left| \frac{\bar{\alpha}_i}{\bar{p}} \right| = \left| \frac{\bar{a}_i \bar{x}^{N-i}}{\bar{p}} \right| = \frac{|a_0| \sum (|z_{j_1}+M| \cdots |z_{j_i}+M|) |x+M|^{N-i}}{|p|}, \quad (9.36)$$

$i = 1, \dots, N$, sekä $|\bar{\alpha}_0/\bar{p}| = |a_0| |x+M|^N / |p|$.

Summan ja sen suurimman jäsenen itseisarvot edustavat yleensä samaa suuruusluokkaa, joten voimme odottaa myös $|\bar{\beta}_i/\bar{p}|$:n ja $|\bar{\gamma}_i/\bar{p}|$:n ja siis kaikkien $\bar{E}_{N_2}^{(i)}$:n lausekkeessa (9.19) esiintyvien kertoimien kasvavan $|M|$:n kasvaessa.

Nollakohtien arvot lähestyvät toisiaan suhteellisesti $|M|$:n suuretessa, mutta myös $|x|$:t kasvavat, joten on syytä tarkastella nollakohtien origosta mitattujen absoluuttisten etäisyyksien vaikutusta kertoimiin.

Ajattelemme kaikkien x -koordinaattien tulevan kerrotuksi kertoimella k y -koordinaattien pysyessä ennallaan. Tällöin polynomia p vastaa polynomi \hat{p} , jolle

$$\hat{p}(\hat{x}) = \hat{p}(kx) = p(x) \quad (9.37)$$

Yhtälön (9.31) mukaan

$$\hat{a}_i = \hat{a}_0 (-1)^i \sum_{j=1}^k (kz_{j_1} \cdots kz_{j_i}) = \frac{\hat{a}_0 k^i}{a_0} a_i, \quad i = 1, \dots, N. \quad (9.38)$$

Tällöin

$$\hat{p}(\hat{x}) = \sum_{i=0}^N \hat{a}_i \hat{x}^{N-i} = \sum_{i=0}^N \frac{\hat{a}_0}{a_0} k^i a_i k^{N-i} x^{N-i} = \frac{\hat{a}_0 k^N}{a_0} p(x), \quad (9.39)$$

joten yhtälöiden (9.37) ja (9.38) perusteella

$$\hat{a}_i = k^{i-N} a_i, \quad i = 0, \dots, N. \quad (9.40)$$

Edelleen

$$\frac{\hat{\alpha}_i}{\hat{p}} = \frac{\hat{a}_i \hat{x}^{N-i}}{\hat{p}} = \frac{k^{i-N} a_i k^{N-i} x^{N-i}}{p} = \frac{\alpha_i}{p} \quad (9.41)$$

ja yhtälöiden (9.10) perusteella myös $\hat{\beta}_i/\hat{p} = \beta_i/p$, ja $\hat{\gamma}_i/\hat{p} = \gamma_i/p$, $i = 0, \dots, N$. Kaikki E :n kertoimet pysyvät siis ennallaan mielivaltaisella k :n arvolla.

Myöskään y -koordinaattien kertominen vakioarvolla x -koordinaattien pysyessä ennallaan ei vaikuta $E_{N,1}^{(0)}$:n kertoimiin, sillä $kp(x) = ka_0(x-z_1) \cdots (x-z_N)$ eli tämä kertominen voidaan samaistaa a_0 :n kertomiseen vakioarvolla. Voimme helposti todeta, että

$$\frac{\alpha_i}{p} = \frac{a_0 (-1)^i \sum (z_1 \cdots z_i) x^{N-i}}{a_0 (x-z_1) \cdots (x-z_N)} \quad (9.42)$$

on riippumaton a_0 :sta.

Edellä saatujen tulosten perusteella $E_{N,2}^{(0)}$:n kertoimet määräytyvät yksinomaan x :n ja nollakohtien arvojen suhteista toisiinsa. ~~ka Nollakohtienäollessa~~ ~~saavsaamanmerkkisiä~~ ~~läheläestyessä~~ ~~ntöisissä~~ ~~kasva-~~ ~~vat~~ ~~kertoimet~~ ~~yleensä~~ ~~itseisarvoltaan.~~

Tarkastellessamme tietyn polynomin p eri pisteitä ovat erityisen mielenkiintoisia polynomin nollakohtien ohella $|x|$:n erittäin suuret arvot sekä piste $x = 0$. Yleisenä huomiona voidaan todeta, että $e_{N,2}$:n kerroin β_N/p on aina yksi.

Jos a_i ei ole nolla, α_i/p ($i = 0, \dots, N$) kasvaa rajatta x :n lähestyessä polynomin nollakohtaa. Samalla yleensä myös β_i/p ja γ_i/p kasvavat rajatta.

Kun $|x|$ kasvaa rajatta, $\alpha_0/p \rightarrow 1$ ja $\alpha_i/p \rightarrow 0$, $i = 1, \dots, N$, joten

$$\lim_{|x| \rightarrow \infty} E_{N,2}^{(1)} = Ne_x + e_0 + \sum_{i=1}^N e_{i,1} + \sum_{i=1}^N e_{i,2} \quad (9.43)$$

ja kertoimet riippuvat siis ainoastaan polynomin asteluvusta.

Jos $a_k \neq 0$ ja $a_j = 0$, $j > k$, niin x :n lähestyessä nollaa $\alpha_k/p \rightarrow 1$ ja $\alpha_j/p \rightarrow 0$, $j \neq k$. Kun otamme huomioon, että tällöin $e_{j,2} = 0$, $j > k$, saamme

$$\lim_{x \rightarrow 0} E_{N,2}^{(1)} = (N-k)e_x + e_k + \sum_{i=k+1}^N e_{i,1}, \quad (9.44)$$

joten kertoimet riippuvat vain asteluvusta ja indeksin k arvosta. Tapaus $k = 0$, jolloin $p = a_0 x^N$, on erityisen huomion arvoinen. Tällöin E :n lauseke on x :n arvosta riippumatta muotoa

$$E_{N,2}^{(1)} = Ne_x + e_0 + \sum_{i=1}^N e_{i,1}. \quad (9.45)$$

Nollakohtien lähekkäisyyden vaikutusta tutkittiin kokeellisesti kolmannen asteen polynomin avulla, jonka alkuperäiset nollakohdat olivat väliltä $(-5,5)$ valitut satunnaiset pisteet

$$z_1 = -4.0467951$$

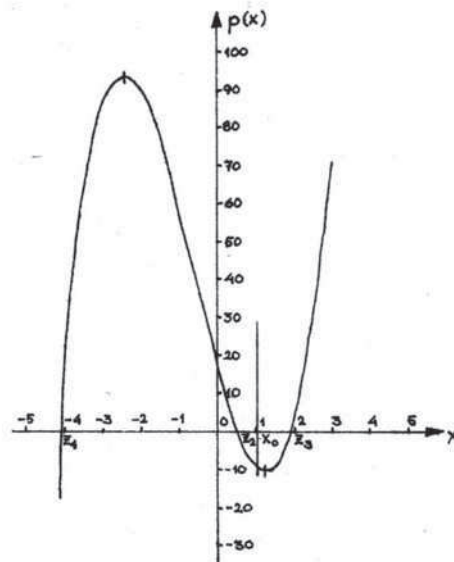
$$z_2 = 0.54756939$$

$$z_3 = 1.9969324$$

ja alkuperäinen a_0 .

$$a_0 = 3.9603187.$$

Alkuperäiselle polynomille (kuva 8) suoritettiin yhtälön (9.32) mukainen siirto, kun M sai arvot 2^m , $m = -7, -6, \dots, 6, 7$. Kullakin m :n arvolla laskettiin polynomin kertoimet a_1, a_2 ja a_3 käyttäen kaavaan (9.31) perustuvaa kertomienlaskualgoritmia K .



kuva 8

Algoritmi K (Polynomin kertoimien lasku). Algoritmi laskee N :nnen asteen polynomin $p = a_0 x^N + a_1 x^{N-1} + \dots + a_{N-1} x + a_N$ kertoimet $a_0 : n$ ja nollakohtien z_1, z_2, \dots, z_N avulla käyttäen apuvektoria s_1, s_2, \dots, s_N .

K1. [Alkuasetukset.] $s_i \leftarrow z_i, i = 1, \dots, N, \text{sign} \leftarrow -1, i \leftarrow 1$.

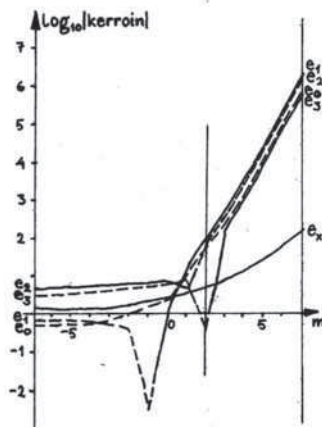
K2. [Kertoimen a_i lasku.] $a_i \leftarrow a_0 * \text{sign} * \sum_{j=1}^{N-i+1} s_j$.

K3. [Seuraava kerroin.] Jos $i = N$, algoritmi päättyy. Muuten $\text{sign} \leftarrow -\text{sign}, i \leftarrow i+1, j \leftarrow 1$.

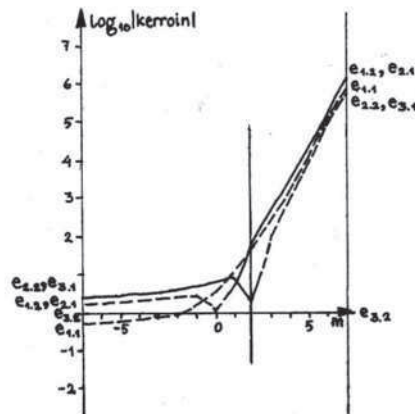
K4. [Apuvektorin täyttö.] $s_j \leftarrow z_j * \sum_{k=j+1}^{N-i+2} s_k$. Jos $j = N-i+1, \rightarrow K2$, muuten $j \leftarrow j+1, \rightarrow K4$. ■

Saaduille polynomeille suoritettiin myös yhtälön (9.37) mukainen supistus siten, että satunnainen piste $x_0 = 1.0768227$ pysyi paikoillaan kaikilla m :n arvoilla. Tällöin k sai kullakin M :n arvolla arvon $x_0 / (x_0 + M)$.

Taylor-sarjan (9.19) kertoimia tutkittiin useilla eri x :n arvoilla, joille polynomin siirtyessä suoritettiin vastaavat siirrot. Suoritettu supistus ei odotusten mukaisesti vaikuttanut kertoimiin.

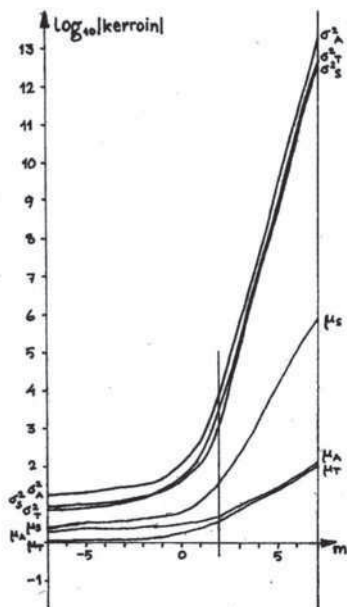


kuva 9a



kuva 9b

Kuvissa 9a ja 9b nähdään kertoimien itseisarvojen logaritmit tekijöittäin pisteessä x_0 . m :n vaihdellessa -7 :stä 7 :ään. Nollakohtien ollessa erimerkkisiä eivät kertoimet muutu oleellisesti. Kaikki nollakohdat tulevat samanmerkkisiksi suunnil-



kuva 10

leen $m:n$ arvolla kaksi, jolloin $M = 4$, ja alkavat tämän jälkeen lähestyä toisiinsa suhteellisesti, jolloin kertoimien itseisarvot alkavat odotusten mukaan kasvaa voimakkaasti. Eräiden kertoimien etumerkki vaihtui tarkastelun kuluessa. Kuvissa 9 merkitsee ehyt viiva positiivista ja katkoviiva negatiivista kerrointa.

Kuvassa 10 on vastaavasti yhtälöiden (9.20) ja (9.21) määrittämien $E_{N,2}^{(1)}$:n odotusarvon ja varianssin lausekkeiden kertoimet, jotka antavat kokonaiskuvan nollakohtien tiivistymisen vaikutuksesta Hermer-shaman kumulatiiviseen pyöristysvirheeseen.

Myös muissa tarkastelluissa pisteissä olivat tulokset vastaavanlaisia. Tosin kertoimet alkoivat kasvaa pienintä nollakohtaa pienemmillä $x:n$ arvoilla vasta näidenkin muuttuessa lisäyksen M vaikutuksesta positiivisiksi. Tämä olikin yhtälön (9.36) mukaan edellytyksenä kertoimien kasvulle.

Vastaava koe suoritettiin myös viidennentoista asteen polynomilla, mutta laskentatarkkuus ei riittänyt luotettavien arvojen saamiseen kertoimille. Esimerkkinä mainittakoon, että nollakohtien ollessa lähimmäimmillään ($m = 7$) saatiin polynomin arvoksi sen eräessä nollakohdassa $4.2 \cdot 10^{28}$. Tämäkin tosin omalla tavallaan osoittaa, että pyöristysvirheet kasvavat voimakkaasti nollakohtien lähestyessä toisiaan.

10. MATRIISIN KÄÄNTÖ

Käännettäessä $m \times m$ -matriisi Gauss-Jordanin menetelmällä [7] otetaan käyttöön apumatriisina $m \times m$ -yksikkömatriisi. Nämä kaksi matriisia yhdessä muodostavat $m \times 2m$ -matriisin

$$(A|I) = \left(\begin{array}{ccc|ccc} a_{11}^{[0]} & a_{12}^{[0]} & \dots & a_{1m}^{[0]} & a_{1,m+1}^{[0]} & \dots & a_{1,2m}^{[0]} \\ a_{21}^{[0]} & a_{22}^{[0]} & \dots & a_{2m}^{[0]} & a_{2,m+1}^{[0]} & \dots & a_{2,2m}^{[0]} \\ \vdots & \vdots & & \vdots & \vdots & & \vdots \\ a_{m1}^{[0]} & a_{m2}^{[0]} & \dots & a_{mm}^{[0]} & a_{m,m+1}^{[0]} & \dots & a_{m,2m}^{[0]} \end{array} \right) \quad (10.1)$$

missä $a_{ij}^{[0]} = a_{ij}$, kun $j \leq m$, $a_{ij}^{[0]} = 1$, kun $i = j - m$ ja $a_{ij}^{[0]} = 0$ muulloin.

Matriisin kääntö tapahtuu algoritmin

$$\begin{cases} a_{ij}^{[k]} = a_{ij}^{[k-1]} - \frac{a_{ik}^{[k-1]} a_{kj}^{[k-1]}}{a_{kk}^{[k-1]}}, & i = 1, \dots, m, i \neq k, \\ & j = k+1, \dots, k+m, \\ & k = 1, \dots, m \\ a_{ij}^{[m+1]} = \frac{a_{ij}^{[m]}}{a_{ii}^{[m]}}, & i = 1, \dots, m \\ & j = m+1, \dots, 2m \end{cases} \quad (10.2)$$

avulla. Mikäli k :nnella iteraatiokierroksella $a_{ij}^{[k]}$:lle ei lasketa uutta arvoa, on $a_{ij}^{[k]} = a_{ij}^{[k-1]}$. Tällöin saamme käänteismatriisiksi

$$A^{-1} = \left(\begin{array}{ccc} a_{1,m+1}^{[m+1]} & a_{1,m+2}^{[m+1]} & \dots & a_{1,2m}^{[m+1]} \\ \vdots & \vdots & & \vdots \\ a_{m,m+1}^{[m+1]} & a_{m,m+2}^{[m+1]} & \dots & a_{m,2m}^{[m+1]} \end{array} \right) \cdot \quad (10.3)$$

Alkioiden indeksointijärjestystä sopivasti vaihtamalla on mahdollista suorittaa kääntäminen tehokkaammin, so. saada tulos tarkemmaksi. Seuraavassa tarkastelussa tätä ei ole otettu huomioon.

Kun hajoitamme algoritmin (10.2) yksittäisiin laskutoimituksiin, saamme sen muotoon

$$\left\{ \begin{array}{l} q_{i,j}^{[0]} = a_{i,j}^{[0]}, \quad i=1, \dots, m, \quad j=1, \dots, 2m \\ q_{i,o,1}^{[k]} = a_{i,k}^{[k-1]} / a_{k,k}^{[k-1]} \\ q_{i,j,2}^{[k]} = q_{i,o,1}^{[k]} \cdot a_{k,j}^{[k-1]} \\ q_{i,j,3}^{[k]} = a_{i,j}^{[k]} - a_{i,j}^{[k-1]} - q_{i,j,2}^{[k]} \\ q_{i,j}^{[m+1]} = a_{i,j}^{[m+1]} = a_{i,j}^{[m]} / a_{i,i}^{[m]} \end{array} \right\} \begin{array}{l} i=1, \dots, m, \quad j=1, \dots, 2m \\ \\ \\ j=k+1, \dots, k+m \\ \\ \\ i=1, \dots, m, \quad i \neq k, \\ k=1, \dots, m \\ \\ i=1, \dots, m, \quad j=m+1, \dots, 2m \end{array} \quad (10.4)$$

jossa q :n indekseistä ensimmäinen osoittaa riviä, toinen saraketta ja kolmas laskutoimituksen vaihetta. Kuhunkin q :hun liittyy täsmälleen yksi yksittäinen pyöristysvirhe. Pyöristysvirheiden lukumääräksi saamme tyypeittäin

tyyppi	lukumäärä
$e_{i,j}^{[0]}$	$m \cdot 2m$
$e_{i,o,1}^{[k]}$	$(m-1) \cdot m$
$e_{i,j,2}^{[k]}$	$(m-1) \cdot m \cdot m$
$e_{i,j,3}^{[k]}$	$(m-1) \cdot m \cdot m$
$e_{i,j}^{[m+1]}$	$m \cdot m$

Kaikkiaan pyöristysvirheitä ja samalla ensimmäisenasteen termejä kumulatiivisen pyöristysvirheen Taylor-kehitelmässä on $2m^3 + 2m^2 - m$ kappaletta.

Taylor-kehitelmän kertoimien analyttinen määrittäminen kullekin käänteismatriisin alkiolle on algoritmin (10.4) laajuuden johdosta varsin hankalaa, mutta kokeellisissa tarkasteluissa voidaan käyttää esimerkiksi kertoimienlaskualgoritmia L. Esimerkiksi 10×10 -matriisin arvojen ollessa valmiina VA-taulukossa saamme alkion 2.7 (E,e)-sarjan taulukkaan C1 sekä alkion 7.2 (E,e)-sarjan taulukkaan C2 ja (R,e)-sarjan taulukkaan C3 ohjelmalla

C MÄÄRITTELYT

REAL VA(10,10), V(2200), C1(2200)

REAL C2(2200), C3(2200)

INTEGER A(10,10), T(2200), G1(2200), G2(2200), Q

C ALKUARVOT

```
I=LBEGIN(V,T,C1,C2,C3,2200)
DO 30 I=1,10
DO 10 J=1,10
10 A(I,J)=LNAME(VA(I,J))
DO 30 J=11,20
VAL=0.
IF(J-I-10) 30,20,30
20 VAL=1.
30 A(I,J)=LNAMEX(VAL)
C ITERAATIOKIERROKSET 1...10
DO 50 K=1,10
DO 50 I=1,10
IF(I-K) 40,50,40
40 Q=LDIV(A(I,K),A(K,K))
DO 50 L=1,10
J=K+L
A(I,J)=LSUB(A(I,J),LMUL(Q,A(K,J)))
50 CONTINUE
C ITERAATIOKIERRÖS 11
DO 60 I=1,10
DO 60 J=11,20
60 A(I,J)=LDIV(A(I,J),A(I,I))
C KERTOIMIEN LASKU
I=LEND(A(2,17))+LREL(C1,2200)
I=LEND(A(7,12))+LREL(C2,2200)+LABS(C3,2200)
```

Yksittäisten pyöristysvirheiden lukumäärä kasvaa niin voimakkaasti matriisin ulottuvuuden kasvessa, että algoritmi L on sovellettava käyttämään tukimuisteja, mikäli sitä aiotaan käyttää suurilla matriiseja käännettäessä.

Algoritmia L hyväksi käyttäen laskettiin (E,e)-sarjat symmetrisen 5x5-matriisin $A = (a_{ij})$ käänteismatriisin alkiuille. Matriisin A satunnaislukualkiot olivat neljän desimaalinumeron tarkkuudella

$$A = \begin{bmatrix} 0.5758 & -0.1035 & -0.0824 & 0.1051 & -0.2077 \\ -0.1035 & 0.2601 & -0.0850 & -0.1125 & 0.2108 \\ -0.0824 & -0.0850 & 0.3935 & 0.0737 & -0.1073 \\ 0.1051 & -0.1125 & 0.0737 & 0.3049 & -0.1401 \\ -0.2077 & 0.2108 & -0.1073 & -0.1401 & 0.3595 \end{bmatrix} .$$

Käänteismatriisiksi $A^{-1} = (c_{ij})$ saatiin

$$A^{-1} = \begin{bmatrix} 2.5555 & -0.2655 & 1.0407 & -0.4111 & 1.7821 \\ -0.2655 & 7.6127 & 0.2779 & 0.9140 & -4.1777 \\ 1.0407 & 0.2779 & 3.2415 & -0.4788 & 1.2193 \\ -0.4111 & 0.9140 & -0.4788 & 4.2070 & 0.7235 \\ 1.7821 & -4.1777 & 1.2193 & 0.7235 & 6.9063 \end{bmatrix} .$$

Saatujen Taylor-kehittelmiin suurimmat kertoimet olivat alkion $c_{4,2}$ sarjassa virheen $e_{5,2}^{[0]}$ (eli alkuarvon $a_{5,2}$ pyöristysvirheen) sekä alkion $c_{2,4}$ sarjassa virheen $e_{2,5}^{[0]}$ (eli alkuarvon $a_{2,5}$ pyöristysvirheen) kertoimet. Näiden molempien arvo oli 10.7718. Noin kaksi kolmasosaa kaikista kertoimista oli nollia.

Kokonaiskuvan saaminen yksittäisten kertoimien perusteella on niiden lukuisuuden vuoksi varsin vaikeaa. Tässä suhteessa saamme paremman käsityksen yhtälöiden (5.8) ja (5.9) antamien keskiarvon ja varianssin kertoimien avulla. Tällöin on syytä huomioida, että mahdollisista virhelähteistä on m^2 kappaletta tarkkoja alkuarvoja (yksikkömatriisin alkiot), $m \cdot (m-1)$ kappaletta muotoa $0-q_{i,j,2}^{[0]}$ olevia ja siten tarkkoja vähennyslaskuja sekä $m \cdot (m-1)$ kappaletta muotoa $q_{i,0,1}^{[0]}$ olevia tarkkoja kertolaskuja. Näitä alkuarvoja ja laskutoimituksia vastaavat yksittäiset pyöristysvirheet ovat nollia, joten niiden kertoimet, yhteensä $3m^2 - 2m$ kappaletta, on syytä poistaa keskiarvon ja varianssin lausekkeita laskettaessa. Tähän on varauduttu aliohjelmaryhmässä L (vastaavat TYPE-kentät ovat negatiivisia).

Esimerkkimatriisille A^{-1} saatiin odotusarvojen lausekkeiden kertoimiksi kahden desimaalin tarkkuudella, kun mainitut poistot oli suoritettu (kertoimet tekijöittäin järjestyksessä $\mu_A/\mu_S/\mu_T$)

-1.00/3.68/2.32	-1.00/2.00/12.68	-1.00/1.00/3.81	-1.00/0.00/3.18	-1.00/-1.00/3.65
-1.00/2.00/12.68	-1.00/1.54/2.95	-1.00/0.00/-3.33	-1.00/-1.00/2.79	-1.00/-2.00/3.67
-1.00/1.00/3.81	-1.00/0.00/-3.33	-1.00/-0.09/1.68	-1.00/-2.00/2.48	-1.00/-3.00/3.36
-1.00/0.00/3.18	-1.00/-1.00/2.79	-1.00/-2.00/2.48	-1.00/-2.02/1.23	-1.00/-4.00/2.23
-1.00/-1.00/3.65	-1.00/-2.00/3.67	-1.00/-3.00/3.36	-1.00/-4.00/-2.23	-1.00/-4.00/2.93

Vastaavat varianssien lausekkeiden kertoimet olivat tekijöittäin järjestyksessä $\sigma_A^2/\sigma_S^2/\sigma_T^2$

2.71/3.34/1.29	366.14/326.72/179.08	5.76/6.51/2.34	22.41/16.71/7.95	14.37/11.31/4.40
366.14/326.72/179.08	6.19/7.53/1.95	207.84/264.30/95.22	32.69/29.60/8.00	18.31/18.07/4.14
5.76/6.51/2.44	207.84/264.30/155.25	1.76/4.74/1.08	12.75/18.96/5.93	16.46/20.74/7.28
22.41/16.71/13.82	32.69/29.60/16.66	12.75/18.96/6.80	17.02/4.62/1.05	56.18/41.51/20.69
14.37/11.31/7.20	18.31/18.07/7.17	16.46/20.74/7.37	56.18/41.51/16.57	8.59/3.65/2.78

Edellä esitetyistä matriiseista voidaan tehdä mielenkiintoisia havaintoja:

- Kaikki kertoimet σ_T^2 :n kertoimia lukuunottamatta ovat symmetrisiä (so. c_{ij} :n tietty kerroin on sama kuin c_{ji} :n vastaava kerroin).
- σ_T^2 :n kertoimet yläkolmiossa ovat pienempiä kuin symmetriset kertoimet alakolmiossa (paitsi c_{45} :n σ_T^2 :n kerroin), joten muiden kertoimien yhtäsuuruudesta johtuen yläkolmioon lasketut käänteismatriisin alkiot ovat tarkempia kuin alakolmioon lasketut.
- Varianssin kerroinmatriisin kerrointen suuruusluokka on päälävistäjällä pienempi kuin muualla, so. päälävistäjän alkioiden arvot ovat tarkempia kuin muiden alkioiden.
- μ_S :n kertoimien arvo alkioidissa c_{ij} , $i \neq j$, on $m-1-j$ ($m = 5$) sekä alkiossa c_{mm} $m-1$.
- Kaikkien μ_A :n kertoimien arvo on -1 .

Jätän tässä esityksessä avoimeksi, ovatko nämä havainnot sattumia, vain käsiteltyyn matriisiin liittyviä, vai voidaan vastaavia havaintoja tehdä yleisesti. Odotusarvomatriisia koskevat havainnot voidaan yleistää ainakin mielivaltaiselle 2×2 -matriisille, sillä tällöin saadaan odotusarvomatriisiksi

$$\left[\begin{array}{cc} -1/\frac{D}{a_{11}a_{22}} / 1 + 2\frac{a_{12}a_{21}}{D} + \frac{a_{12}a_{21}}{a_{11}a_{22}} & -1/ -1 / 2\frac{a_{11}a_{22}}{D} \\ -1/ -1 / 2\frac{a_{11}a_{22}}{D} & -1/ +1 / 1 + 2\frac{a_{12}a_{21}}{D} \end{array} \right]$$

missä $D = a_{11}a_{22} - a_{12}a_{21}$ on ko. matriisin determinantti.

Liite: ALGORITMI L FORTRAN IV-OHJELMANA

algoritmin es
askel

```
L1  C      FUNCTION SUBPROGRAM GROUP L
    C
    C      INTEGER FUNCTION QI
    C      ENTRY LBEGIN (VALUE,TYPE,OPER1,OPER2,COEFF,M)
    C      INTEGER TYPE,OPER1,OPER2,QJ,QK,QN,EX
    C      DIMENSION VALUE(M),TYPE(M),OPER1(M),OPER2(M)
    C      DIMENSION COEFF(M)
    C      EX=0
    C      QI=0
    C      I=1
    C      RETURN

L3  C      INITIAL VALUES
    C
    C      ENTRY LNAME(VAL)
    C      TYPE(I)=1
    C      10 VALUE(I)=VAL
    C      GO TO 66
    C      ENTRY LNAMEX(VAL)
    C      TYPE(I)=-1
    C      GO TO 10

L4  C      NEGATION
    C
    C      ENTRY LNEG(QJ)
    C      TYPE(I)=-2
    C      VALUE(I)=-VALUE(QJ)
    C      GO TO 65

L5  C      ADDITION
    C
    C      ENTRY LADD(QJ,QK)
    C      IF(VALUE(QJ)) 20,21,20
    C      20 IF(VALUE(QK)) 22,21,22
    C      21 EX=1
    C      22 TYPE(I)=3
    C      VALUE(I)=VALUE(QJ)+VALUE(QK)
    C      GO TO 60
    C      ENTRY LADDX(QJ,QK)
    C      GO TO 21
    C      ENTRY LADDN(QJ,QK)
    C      EX=-1
    C      GO TO 22
```

```
L6      C
        C      SUBTRACTION
        C
        ENTRY LSUB(QJ,QK)
        IF(VALUE(QJ)) 30,31,30
30      IF(VALUE(QK)) 32,31,32
31      EX=1
32      TYPE(I)=4
        VALUE(I)=VALUE(QJ)-VALUE(QK)
        GO TO 60
        ENTRY LSUBX(QJ,QK)
        GO TO 31
        ENTRY LSUBN(QJ,QK)
        EX=-1
        GO TO 32

L7      C
        C      MULTIPLICATION
        C
        ENTRY LMUL(QJ,QK)
        IF(ABS(VALUE(QJ))-1.) 40,41,40
40      IF(ABS(VALUE(QK))-1.) 42,41,42
41      EX=1
42      TYPE(I)=5
        VALUE(I)=VALUE(QJ)*VALUE(QK)
        GO TO 60
        ENTRY LMULX(QJ,QK)
        GO TO 41
        ENTRY LMULN(QJ,QK)
        EX=-1
        GO TO 42

L8      C
        C      DIVISION
        C
        ENTRY LDIV(QJ,QK)
        IF(ABS(VALUE(QK))-1.) 51,50,51
50      EX=1
51      TYPE(I)=6
        VALUE(I)=VALUE(QJ)/VALUE(QK)
        GO TO 60
        ENTRY LDIVX(QJ,QK)
        GO TO 50
        ENTRY LDIVN(QJ,QK)
        EX=-1
        GO TO 51

L9      C
        60 OPER2(I)=QK
        IF(VALUE(I)) 62,61,62
        61 EX=EX+1
        62 IF(EX) 64,64,63
        63 TYPE(I)=-TYPE(I)
        64 EX=0
L10     65 OPER1(I)=QJ
L11     66 QI=I
        I=I+1
        RETURN
```

```

C
C
C      COEFFICIENTS

L12      ENTRY LEND(QN)
          DO 80 K=1,M
          80 COEFF(K)=0.
            COEFF(QN)=1.
            K=QN

L13      81 ITYP=IABS(TYPE(K))
          K1=OPER1(K)
          K2=OPER2(K)

L14      GO TO (83,102,103,104,105,106),ITYP
L15      102 COEFF(K1)=COEFF(K1)-COEFF(K)
L16      103 COEFF(K)=0.
          GO TO 83

L17      104 D1=1.
          D2=1.
          GO TO 82

L18      105 D1=1.
          D2=-1.
          GO TO 82

L19      106 D1=1./VALUE(K2)
          D2=-VALUE(K1)/VALUE(K2)**2

L20      82 COEFF(K1)=COEFF(K1)+D1*COEFF(K)
          COEFF(K2)=COEFF(K2)+D2*COEFF(K)

L21      83 K=K-1
          IF(K) 84,84,81
          84 I=1
            QI=0
            RETURN

C
C
C      RELATIVE OR ABSOLUTE COEFFICIENTS

L22      ENTRY LREL(COEFIC,M)
          RES=VALUE(QN)
          ASSIGN 91 TO IG
          GO TO 90
          ENTRY LABS(COEFIC,M)
          ASSIGN 92 TO IG
          90 CONTINUE
            DIMENSION COEFIC(M)
            DO 92 K=1,QN
              COEFIC(K)=COEFF(K)*VALUE(K)
              GO TO IG,(91,92)
            91 COEFIC(K)=COEFIC(K)/RES
            92 CONTINUE
              QI=0
              RETURN
            END
```

KÄYTETYJÄ MERKINTÖJÄ

merkintä	selitys	sivu
*	pyöristyssymboli, esim. $z^* = z(1+e)$	3
< >	sarjan jäännöstermi, sulkeissa sarjan muuttujatyyppi ja pienin jäännöstermissä esiintyvä asteluku, esim. $\langle e^3 \rangle$	15
-	x-koordinaattien siirto: $\bar{x} = x+M$	46
^	x-koordinaattien lavennus/supistus: $\hat{x} = kx$	47
α	polynomille p $\alpha_i = x^{n-i} a_i$	40
β	polynomille p $\beta_i = \int_{j=0}^i x^{n-i} a_j$	40
γ	polynomille p $\gamma_i = \sum_{j=0}^{N-j} x^{n-i} a_j$	40
ε	mantissan absoluuttinen pyöristysvirhe pyöristävässä aritmetiikassa	4
ε'	mantissan absoluuttinen pyöristysvirhe katkaisevassa aritmetiikassa	5
μ	odotusarvo; μ_λ :alkuarvon, μ_s :summan ja erotuksen, μ_T :tulon ja osamäärän pyöristysvirheen odotusarvo	8,9,11, 12,16,18
σ	keskihajonta; σ^2 :varianssi; σ_λ^2 :alkuarvon, σ_s^2 :summan ja erotuksen, σ_T^2 :tulon ja osamäärän pyör.virheen varianssi	8,9,11, 12,16,18
ξ, ζ, η	satunnaismuuttujia	7,8,16,45
a	polynomien p kerroin; matriisin A alkio; (E,e)-sarjan kerroin	39,52 15
b	käytetyn aritmetiikan kantaluku	2
c	(R,e)-sarjan kerroin; matriisin A ⁻¹ alkio	20,55
D	Q:n derivaatta operandinsa suhteen, esim. $D = \partial Q_n(q_i, q_j) / \partial q_i$;	20
δ	keskihajonta; D^2 : varianssi, esim. $D^2(e)$	8
\bar{d}	(R,r)-sarjan kerroin	22
E	kumulatiivinen suhteellinen pyör.virhe; odotusarvo, esim. E(e)	15 8
$E_n^{(i)}$	(E,e)-sarjan i:nneen asteen termien summa	15

(E,e)	(E,e)-sarja: muotoa $E_n = \sum_{t=0}^n a_t e_t + \dots$ oleva Taylor-kehitemä	15
e	yksittäinen suhteellinen pyör.virhe, esim. $z^* = z(1+e)$	3
f	tiheysfunktio	5
i,j,k,l,m,n,N	indeksejä	
k	x-koordinaattien lavennus/supistusker- roin, vrt. " \sim "	47
M	x-koordinaattien siirron määrä, vrt. " \sim "	46
m	liukuluvun mantissa	2
p	liukuluvun eksponentti; polynomi $p = p(x) = a_0 x^n + a_1 x^{n-1} + \dots + a_n$	2 39
Q	laskutoimitus, esim. voi olla $Q(a,b)=a-b$	19
q	algoritmin alkuarvo tai tulos	19
R	kumulatiivinen absoluuttinen pyör.virhe	19
$R_n^{(i)}$	(R,e)-sarjan i:nneen asteen termien summa	20
(R,e)	(R,e)-sarja: muotoa $R_n = \sum_{t=0}^n c_t e_t + \dots$ oleva Taylor-kehitemä	20
(R,r)	(R,r)-sarja: muotoa $R = \sum_{t=0}^n d_t r_t + \dots$ oleva Taylor-kehitemä	22
r	yksittäinen absoluuttinen pyör.virhe	3
t	liukuluvun mantissan numeroiden lkm.	2
u	b^{-t}	4
x	piste, jossa polynomin p arvo lasketaan	39
z	liukuluku, $z = m \cdot b^p$; polynomin nollakohta	2 46

VIITELUETTELO

- [1] Babuška, I., Numerical Stability in Mathematical Analysis, Information Processing 68, 11-23, Amsterdam, 1969
- [2] Elfving, G., Todennäköisyyslaskenta, II luku, Otava, Helsinki, 1966
- [3] Hamming, Numerical Methods for Scientists and Engineers, luku 2, McGraw-Hill, New York, 1962
- [4] Henrici, P., Elements of Numerical Analysis, luvut 15 ja 16, Wiley, New York, 1964
- [5] Henrici, P., Discrete Variable Methods in Ordinary Differential Equations, Wiley, New York, 1962
- [6] Hull, T. E. ja Swenson, J. R., Tests of Probabilistic Models for Propagation of Round-off Errors, Comm. ACM, vol.9, 108-113, 1966
- [7] Isaacsson, E. ja Keller, H. B., Analysis of Numerical Methods, luvut 1 ja 2, Wiley, New York, 1966
- [8] Knuth, D. E., The Art of Computer Programming, vol.2, luku 4, Addison-Wesley, New York, 1969
- [9] Tienari, M., A Statistical Model of Roundoff Errors in Varying Length Floating-point Arithmetic, Helsinki, 1970
- [10] Wilkinson, J., Rounding Errors in Algebraic Processes, Prentice Hall, Englewood Cliffs, N.J., 1963