

บทที่ 1
คอมพิวเตอร์ ความคลาดเคลื่อน
และขั้นตอนวิธี

- 1.1 วิธีวิเคราะห์เชิงตัวเลข (Numerical Method) คืออะไร ?
- 1.2 เครื่องคำนวณแบบดิจิทัล (digital devices) ทำผิดหรือไม่ ?
- 1.3 เครื่องคอมพิวเตอร์เก็บเลขจำนวนจริง (real numbers) อย่างไร ?
- 1.4 ความคลาดเคลื่อนในการคำนวณแบบ Fixed-Precision
- 1.5 การหลีกเลี่ยงความคลาดเคลื่อนจากการปัดเศษ (Roundoff Error)
- 1.6 ความคลาดเคลื่อน ความแม่นยำ และ การทดสอบเพื่อตรวจความใกล้เคียง
- 1.7 ขั้นตอนวิธีของการทำซ้ำ (Iterative Algorithms)

บทที่ 1

คอมพิวเตอร์ ความคลาดเคลื่อน และขั้นตอนวิธี

1.1 วิธีวิเคราะห์เชิงตัวเลข (Numerical Method) คืออะไร ?

ในวิชาแคลคูลัสเบื้องต้นเราเรียนวิธีการหาอนุพันธ์และการอินทิเกรตเพื่อให้ได้ผลเฉลยแน่นอนตรง (exact answers) แต่โชคไม่ดีที่เทคนิคเบื้องต้นทางแคลคูลัส (และไม่ว่าจะขั้นสูง) เพียงอย่างเดียวไม่พอเพียงสำหรับการแก้ปัญหาแคลคูลัสตามตัวอย่างต่อไปนี้

ปัญหาที่ 1 หา maximum และ minimum values ซึ่งอยู่ในช่วง $[0,1]$ ของ

$$F(x) = x^6 + 5x^4 - 9x + 1$$

ในการทำนี้เราต้องแก้สมการ polynomial degree 5 นั่นคือ

$$dF(x)/dx = F'(x) = 6x^5 + 20x^3 - 9$$

ปัญหาที่ 2 หาค่า $\int_0^b e^{-x} dx$ และ $\int_0^b \sin(x)/x dx$ สำหรับค่า $b > 0$

ปัญหาที่ 3 หา $y = y(t)$ ซึ่งสอดคล้องกับสมการเชิงอนุพันธ์

$$dy/dx = y^2 + t^2, \quad y(0) = 1 \quad (\text{i.e., } y = 1 \text{ at } t = 0)$$

ปัญหาที่ 4 หา local extrema ของฟังก์ชันของตัวแปร 2 ตัว

$$f(x,y) = xe^y - \ln(2x^2 + y^2) \quad [\ln = \log_e]$$

ซึ่งต้องทำโดยการหารากของสมการ 2 สมการต่อไปนี้คือ

$$\partial f / \partial x = e^y - [4x / (2x^2 + y^2)] = 0$$

$$\partial f / \partial y = xe^y - [2y / (2x^2 + y^2)] = 0$$

ถ้าเราลองแก้ปัญหาคำที่ 1-4 โดยวิธีที่เรียนมา จะพบว่าเป็นไปได้ที่จะหา

สูตรที่ชัดเจนซึ่งจะใช้กับแต่ละปัญหาเพื่อให้ได้คำตอบที่ถูกต้อง

ในสถานการณ์จริงๆนั้น เรามักไม่ต้องการผลเฉลยแม่นยำตรง แต่ต้องการผลเฉลยซึ่งมีความแม่นยำถึงระดับที่ต้องการเท่านั้น เราจึงใช้วิธีการวิเคราะห์เชิงตัวเลขช่วยในการหาคำตอบของปัญหาเช่นดังกล่าวข้างต้นได้

1.1A วิธีการวิเคราะห์เชิงตัวเลข และการทดสอบเพื่อหยุดการทำงาน (Termination Tests)

คำ Algorithm ถูกใช้เพื่ออธิบายวิธีการเป็นขั้นๆ (step-by-step) ซึ่งต้องใช้จำนวนขั้นซึ่งมีจำนวนจำกัด (finite number of step)

วิธีการวิเคราะห์เชิงตัวเลข (Numerical method) คือ ขั้นตอนวิธี (Algorithm) หนึ่งสำหรับการหาค่าเชิงตัวเลขซึ่งมีระดับความแม่นยำตามระดับที่ต้องการ

- ตัวอย่าง a. จงอธิบายวิธีการวิเคราะห์เชิงตัวเลขสำหรับการหาอนุพันธ์ของฟังก์ชัน f ณ จุด x โดยให้มีความแม่นยำถึงทศนิยมตำแหน่งที่ 4 (four-decimal-place accuracy)
- b. จงใช้วิธีการวิเคราะห์เชิงตัวเลขสำหรับการประมาณ $d \sin x/dx \Big|_{x=0}$ โดยให้มีความแม่นยำถึงทศนิยมตำแหน่งที่ 4

* เมื่อกล่าวว่า " X ประมาณค่า x ด้วยความแม่นยำถึงทศนิยมตำแหน่งที่ d "

(X approximates x to d decimal places) หมายความว่า X และ x ต่างกันน้อยกว่า $1/2$ unit. ในทศนิยมตำแหน่งที่ d

$$\text{นั่นคือ } |X-x| < (1/2) \cdot 10^{-d}$$

ตัวอย่างเช่น $x = 12.45672$, $X = 12.45676$

$$|X-x| = .00004 < .5 * .0001 = .00005$$

แสดงว่า X ประมาณค่า x ด้วยความแม่นยำถึงทศนิยมตำแหน่งที่ 4

Solution

๑. เมื่อเราไม่มีสูตรทั่วไปสำหรับการหาอนุพันธ์ เราจะพิจารณา
ค่าจำกัดความของ

$$f'(x) = df(x)/dx = \lim_{h \rightarrow 0} \Delta f(x)/h$$

(if it exists)

โดยที่ $\Delta f(x)/h = [f(x+h)-f(x)]/h$ (difference quotient
of f at x)

วิธีแก้ปัญหานี้แบบปกติคือ ประมาณ $f'(x)$ ด้วย $\Delta f(x)/h$ สำหรับ h เล็กๆ
ค่าของ h ที่เล็กพอจะทำให้แน่ใจว่า $\Delta f(x)/h$ ประมาณค่า $f'(x)$ ถึงทศนิยม
ตำแหน่งที่ 4 นั้นขึ้นอยู่กับฟังก์ชัน ที่กำหนดให้ และจุด x ที่ระบุ เพื่อเป็น
การช่วยกำหนดค่าที่เหมาะสมของ h เราเพียงแต่หาค่า $\Delta f(x)/h$ สำหรับค่า h
ที่ลดลงเป็นลำดับ และหยุดเมื่อ $E(h) = f'(x) - \Delta f(x)/h$ มีค่าเข้าใกล้ 0
[$E(h)$ คือความคลาดเคลื่อนของการประมาณ]

ถึงแม้วิธีการข้างต้นจะดูง่ายๆ แต่เราต้องทราบค่าของ $f'(x)$ ซึ่งจริงๆแล้ว
เราไม่ทราบค่าเพราะเรากำลังจะหาค่าของ $f'(x)$ และถ้าเราทราบเราก็ไม่ต้อง
ใช้วิธีวิเคราะห์เชิงตัวเลขเข้าช่วย ดังนั้นเราจะทำการหาค่า $\Delta f(x)/h$ ต่อไป
เรื่อยๆ จนกว่าค่านี้จะ "close enough"

เนื่องจากว่าเราต้องคำนวณ $\Delta f(x)/h$ ซึ่งสำหรับค่าต่างๆของ h 's จะได้
แสดงวิธีดังกล่าวในโปรแกรม DERIV

ในการใช้ DERIV เราต้องระบุ initial h , a shrinking ratio r
และ จำนวนครั้งที่ต้องการทำซ้ำ ของมันต่อไป

```

        Calculate and print h and  $\Delta f(x)/h$ 
        h <-- h/r { replace h by 'current h)/r }
C*****DERIV: A FORTRAN PROGRAM FOR ESTIMATING F'(X)
C   APPROXIMATE F'(X) BY DQ(H), H=HO,HO/R,...,
c   HO/R**(NREPS-1)
        F(X)=SIN(X)          <-- Statement function
        WRITE(6,1)
1   FORMAT(1X,'ENTER X,HO,SHRINKING RATIO.8 REPETITIONS')
        READ(5,*)X,H,R,NREPS
        WRITE(6,2)
2   FORMAT(1X,'      H',11X,'DQ(H)')
        DO 10 N =1,NREPS
        DQ=(F(X+H)-F(X))/H
        WRITE(6,3)H,DQ
3   FORMAT(E13.5,2X,F11.6)
        H=H/R
10  CONTINUE
        STOP
        END

```

Solution

- b. เรามพบปัญหาในการหาค่า h ซึ่งต้องเล็กพอที่จะทำให้ $\Delta f(x)/h$ ประมาณ $f'(x)$ ได้ตามความแม่นยำที่ต้องการ

เมื่อ $f(x) = \sin x$ และ $x = 0$

ถ้าเราเลือก $H_{(initial)} = 1$

R : shrinking ratio 4

และ NREPS = 8

run DERIV ด้วยค่าข้างต้นได้ผลดังนี้

ENTER X,HO,SHRINKING RATIO,# REFETTITIONS

0,1,4,8 ^h column เองเพราะทราบค่า F'(X)=1

H	DQ(H)	E(H)=F'(X)-DQ(H)
0.10000E+01	0.841471	0.156529
0.25000E+00	0.989616	0.010384
1/16 = 0.62500E-01	0.999349	0.000651 = E(1/16)
1/64 = 0.15625E-01	0.999959	0.000041 = E(1/64)
1/256 = 0.39063E-02	0.999997	0.000003
0.97656E-03	1.000000	0.000000
0.24414E-03	1.000000	0.000000
0.61035E-04	1 .000000	0.000000

$h = 1/16 = 0.0625 \rightarrow DQ(h) = 0.999349$ (เอาทศนิยม 4 ตำแหน่ง)

$h = 1/64 = 0.015625 \rightarrow DQ(h) = 0.999959$ (เอาทศนิยม 4 ตำแหน่ง)

ดังนั้น $h = 1/64$ จึงเล็กพอที่จะทำให้แน่ใจว่า $\Delta f(x)/h$ ประมาณ $\sin'(0)$ ได้แม่นยำถึงทศนิยมตำแหน่งที่ 4

(*) โปรดสังเกตว่า

$$DQ(1/64) - DQ(1/16) = 0.999959 - 0.999349 = 0.000610 \approx E(1/16)$$

$$DQ(1/256) - DQ(1/64) = 0.999997 - 0.999959 = 0.000038 \approx E(1/64)$$

โดยทั่วไป ผลต่างระหว่างค่าประมาณที่ต่อเนื่องกัน

นั่นคือ $\Delta DQ = DQ(h) - DQ(\text{preceding } h)$

คือเลขจำนวนซึ่งสามารถหาได้โดยไม่ทราบ $f'(x)$ และจะถูกใช้เพื่อประมาณค่า error of DQ (preceding h) สิ่งที่จะตามมาคือ แทนที่จะต้องปฏิบัติซ้ำทั้ง NREPS ครั้ง เราจะหยุดกระทำเมื่อขนาดของ ΔDQ มีค่าเล็กมาก

(**) เงื่อนไขที่เหมาะสมสำหรับ four-place accuracy คือ

$$|\Delta DQ| < 10^{-4}$$

เมื่อพิจารณา (*) แล้ว จะพบว่า termination test (**) จะเป็นเหตุให้หยุดการกระทำหลังจากพิมพ์ $H = 1/256$ และ $DQ = 0.999997$ แล้ว

1.1B ทำไมต้องใช้ Programmable Device สำหรับการวิเคราะห์เชิงตัวเลข ?

จากตัวอย่างที่ผ่านมาได้แสดงให้เห็นว่าลักษณะ 3 ประการของคอมพิวเตอร์ (หรือ programmable calculator) ดังต่อไปนี้เหมาะสมกับวิธีวิเคราะห์เชิงตัวเลข คือ

1. ความยืดหยุ่น (Flexibility)

โปรดสังเกตว่าเราสามารถเปลี่ยน $F(X)$ ในบรรทัดที่ 4 ของโปรแกรม DERIV ได้ เมื่อโปรแกรมนี้ถูกเก็บในหน่วยความจำ เราประมาณ $f'(x)$ ที่ค่าของ x ใดๆ ด้วยความแม่นยำถึงทศนิยมตำแหน่งที่ต้องการได้ง่าย ๆ โดยการเปลี่ยน $F(X)$ ตามต้องการ แล้ว enter ค่าต่างๆของ H , R และ NREPS ระหว่างการปฏิบัติงาน

2. ความมีประสิทธิภาพ (Efficiency)

ขั้นตอนวิธีใน DERIV เป็นแบบอย่างของขั้นตอนของวิธีการวิเคราะห์เชิงตัวเลข ในแง่ที่ว่ากระบวนการหนึ่ง (calculate the next DQ, output H and DQ แล้ว shrink H) ถูกกระทำจนกระทั่งนิพจน์ที่ใช้เป็นเงื่อนไขของการหยุดทำงานเป็นจริง ขั้นตอนวิธีดังกล่าวถูกเรียกว่า ขั้นตอนวิธีสำหรับการทำซ้ำ ("iterative" algorithm) ในการเขียนโปรแกรมอย่างมีประสิทธิภาพ นั้นทำได้โดยใช้ loop ดังในตัวอย่าง

โปรแกรม DERIV ตั้งแต่

DO 10 N = 1, NREPS

:

10 CONTINUE

3. ความเร็ว และ ความเชื่อมั่น (Speed and Reliability)

หลังจากที่โปรแกรมถูกแก้ไขแล้วและให้ข้อมูลเข้าที่จำเป็น โปรแกรมจะถูกปฏิบัติตามด้วยความเร็วที่เร็วกว่าและเชื่อถือได้มากกว่าที่คนคนใดคนหนึ่งจะทำได้

1.2 เครื่องคำนวณแบบดิจิทัล (Digital Devices) ทำผิดหรือไม่ ?

Digital devices often give erroneous answers.

ข้อครหาหรือเกือบทั้งหมดของ ค่าตอบที่ไม่ถูกต้อง เกิดจาก human error เช่น

เกิดจาก Programmer กำหนดตรรกะ (logic) ผิดในโปรแกรม

เกิดจาก Keypuncher พิมพ์โปรแกรมผิดหรือข้อมูลเข้าผิด

เกิดจาก User ให้ข้อมูลเข้าผิดๆ ใช้สูตรหรือวิธีการที่ผิด

ความคลาดเคลื่อนทั้งหลายนั้นแก้ไขได้โดยการค้นหาให้พบแล้วแก้ไขให้ถูกต้อง

1.2A ตัวอย่างของความคลาดเคลื่อนซึ่งเกิดจากเครื่องคำนวณแบบดิจิทัล

ตัวอย่าง

$$\text{ให้ } A = (1 - \cos x) / x^2, \quad B = \sin^2 x / [x^2 (1 + \cos x)]$$

ถ้าคูณ A ด้วย $(1 + \cos x) / (1 + \cos x)$ จะเห็นว่า A และ B เท่ากัน สำหรับ x ที่ไม่เป็น 0 ทั้งนี้ $\cos x \neq -1$

If $x = 0 + 0.1$, then $A = B = 0.499583$; if $x = \pi + 0.1$, then $A = B = 0.189857$
 If $x = 0 + 0.01$, then $A = B = 0.499996$; if $x = \pi + 0.01$, then $A = B = 0.201353$
 If $x = 0 + 0.001$, then $A = B = 0.500000$; if $x = \pi + 0.001$, then $A = B = 0.202513$
 If $x = 0 + 0.0001$, then $A = B = 0.500000$; if $x = \pi + 0.0001$, then $A = B = 0.202629$
 If $x = 0 + 0.00001$, then $A = B = 0.500000$; if $x = \pi + 0.00001$, then $A = B = 0.202641$

ค่า A และ B ที่แสดงไว้ ได้ถูกปัดเศษให้เหลือทศนิยมเพียง 6 ตำแหน่ง

ต่อไปนี้เป็นตัวอย่างโปรแกรมที่ทำการหาค่าของ A และ B จากค่าของ x ทั้ง

10 คำข้างต้น

```

00100 C * * * * * ERROR * * * * *
00200 C THIS PROGRAM DEMONSTRATES DIGITAL DEVICE ERRORS
00300 c
00400 DATA IW, IR /5, 5/
00500 WRITE(IW,1)
00600 1 FORMAT('ENTER THE LIMITING VALUE OF X')
00700 READ(IR,*) XLIMIT
00800 C
00900 H = 0.1
01000 DO 10 K=1, 5
01100 X = XLIMIT + H
01200 C
01300 A = (1 - COS(X))/X**2
01400 B = (SIN(X)/X)**2/(1 + COS(X))
01500 C
01600 WRITE (IW,2) X, A, B
01700 2 FORMAT(' IF X=',F9.6,' THEN A=',E13.6,
01800 & ' AND B=',E13.6)
01900 H = H/10
02000 10 CONTINUE
02100 STOP
02200 END

ENTER THE LIMITING VALUE OF X
0

IF X= 0.100000 THEN A= 0.499584E+00 AND B= 0.499583E+00
IF X= 0.010000 THEN A= 0.500008E+00 AND B= 0.499996E+00
IF X= 0.001000 THEN A= 0.500639E+00 AND B= 0.500000E+00
IF X= 0.000100 THEN A= 0.745058E+00 AND B= 0.500000E+00
IF X= 0.000010 THEN A= 0.000000E+00 AND B= 0.500000E+00

ENTER THE LIMITING VALUE OF X
3.14159265

IF X= 3.241593 THEN A= 0.189857E+00 AND B= 0.189857E+00
IF X= 3.151593 THEN A= 0.201353E+00 AND B= 0.201346E+00
IF X= 3.142593 THEN A= 0.202513E+00 AND B= 0.199852E+00
IF X= 3.141693 THEN A= 0.202629E+00 AND B= 0.135837E+00
%FRSAPR Floating divide check PC= 235

IF X= 3.141603 THEN A= 0.202641E+00 AND B= 0.170141E+39

```

รูป 1.2 - 1 โปรแกรมภาษาฟอร์แทรนเพื่อหาค่า A และ B ใกล้ 0 และ π

จากตัวอย่างการวิ่งโปรแกรมจะเห็นว่า คอมพิวเตอร์คำนวณค่าของ B ได้แม่นยำเมื่อ x มีค่าใกล้ 0 และคำนวณค่า A ได้แม่นยำเมื่อ x มีค่าใกล้ π และให้ค่าที่ใช้ไม่ได้เมื่อคำนวณค่า A จากค่า x ใกล้ 0 และคำนวณ B จากค่า x ใกล้ π

ดังนั้นจากตัวอย่างแสดงว่า

Computers make no judgment about the correctness of the values they produce; they will display nonsensical answers as neatly and authoritatively as correct ones.

เราจะประสบปัญหา ถ้าเราเชื่อว่า $A = 0$ เมื่อ $x = 0.00001$ เพราะคอมพิวเตอร์ให้คำตอบเช่นนั้น

1.2B ความคลาดเคลื่อนในการใช้สูตรกำลังสอง (Quadratic Formula)

ในการหาราก 2 รากของสมการกำลังสอง $ax^2 + bx + c = 0$ ($a \neq 0$) โปรแกรมและตัวอย่างการวิ่งโปรแกรมต่อไปนี้จะแสดงความคลาดเคลื่อนที่อาจเกิดขึ้นเมื่อหารากของสมการข้างต้น

```

00100 C * * * * * QUAD * * * * *
00200 C THIS PROGRAM DEMONSTRATES ERRORS THAT CAN OCCUR WHEN USING
00300 C THE QUADRATIC FORMULA TO FIND THE REAL ROOTS OF A**2+B**X+C
00400 C
00500 DATA IW, IR /5, 5/
00600 WRITE(IW,1)
00700 1 FORMAT('FOR REAL ROOTS OF AX**2+BX+C, INPUT A, B, C')
00800 READ(IR,*) A, B, C
00900 C
01000 EI = -0.5*B/A
01100 CI = CIA
01200 DISC = B1*B1 - CI
01300 IF (DISC .GE. 0.) GOTO 20
01400 WRITE(IW,2)
01500 2 FORMAT(' ROOTS ARE COMPLEX')
01600 STOP
01700 C
01800 C REAL ROOTS: FIND THEM AS R(+), C1/R(+) AND AS R(-), C1/R(-)
01900 20 D1 = SQRT(DISC)
02000 RPLUS = B1 + D1
02100 RMINU = B1 - D1
02200 R2PLUS = CI/RPLUS
02300 R2MINU = C1/RMINU
02400 C
02500 WRITE(IW,3) RPLUS, RMINU, R2PLUS, R2MINU
02600 3 FORMAT(' ROOTS: R(+) = ',E13.6,5X,' R(-) = ',E13.6,
02700 & /' C1/R(+) = ',E13.6,5X,' C1/R(-) = ',E13.6)
02800 STOP
02900 END

```

รูป 1.2-2 โปรแกรมภาษาฟอร์แทรนเพื่อคำนวณ $r(+)$, $r(-)$, $c_1/r(+)$, และ $c_1/r(-)$

FOR REAL ROOTS OF AX^2+BX+C , INPUT A, B, C
2, 9, -5

ROOTS: $R(+)$ = 0.500000E+00 $R(-)$ = -0.500000E+01
 $C1/R(+)$ = -0.500000E+01 $C1/R(-)$ = 0.500000E+00

FOR REAL ROOTS OF AX^2+BX+C , INPUT A, B, C
1, -100.01, ~

ROOTS: $R(+)$ = 0.100000E+03 $R(-)$ = 0.100002E-01
 $C1/R(+)$ = 0.100000E-01 $C1/R(-)$ = 0.999977E+02

FOR REAL ROOTS OF AX^2+BX+C , INPUT A, B, C
1, +100.01, 1

ROOTS: $R(+)$ = -0.100002E-01 $R(-)$ = -0.100000E+03
 $C1/R(+)$ = -0.999977E+02 $C1/R(-)$ = -0.100000E-01

ในที่นี้ $b_1 = -b/2a$, $c_1 = c/a$ และ $d_1 = \sqrt{b_1^2 - c_1}$
 $r(+)$ = $b_1 + d_1$, $r(-)$ = $b_1 - d_1$
 $r(+)r(-)$ = $b_1^2 - d_1^2 = c_1$
 ดังนั้น $r(+)$ = $c_1/r(-)$, $r(-)$ = $c_1/r(+)$ ทั้งนี้ตัวหารต้อง $\neq 0$

จากตัวอย่างที่ผ่านมาใน 1.2A และ 1.2B ความคลาดเคลื่อนที่เกิดขึ้นจะไม่เป็น
 เอกลักษณะ (unique) สำหรับคอมพิวเตอร์ทุกเครื่อง
 เราต้องทราบที่มาของความคลาดเคลื่อนแล้วแก้ไขหรือทำให้เกิดน้อยที่สุด และเรา
 ต้องรู้วิธีการเก็บเลขจำนวนจริงในหน่วยความจำของคอมพิวเตอร์ รวมทั้งการ
 คำนวณโคสใช้คอมพิวเตอร์ที่ทำการคำนวณด้วย Fixed-Precision

1.3 เครื่องคอมพิวเตอร์เก็บเลขจำนวนจริง (Real Numbers) อย่างไร?

1.3A เลขฐานสอง (Binary Numbers)

เลขฐานสองจำนวนหนึ่งคือเลขจำนวนจริง X ซึ่งอาจถูกเขียนได้ในรูป

$$X = \pm [(b_n 2^n + b_{n-1} 2^{n-1} + \dots + b_1 2^1 + b_0)] + [(b_{-1} 2^{-1} + b_{-2} 2^{-2} + \dots + b_{-k} 2^{-k})]$$

integer part ของ X fractional part ของ X

.....(*)

โดยที่ b_i คือ bit (binary digit) ซึ่งอาจเป็น 0 หรือ 1

นิพจน์ (*) เราเรียกว่า binary expansion of X หรือเรียกย่อ ๆ ว่า

binary representat ion

ดังนั้น

$$X = \pm (b_n b_{n-1} \dots b_1 \overset{\text{binary point}}{b_0} b_{-1} b_{-2} \dots b_{-k})_2 \quad \text{ทุก } b_i \text{ มีค่า 0 หรือ 1}$$

(n+1)-bit integer k-bit fraction

ตัวอย่าง

$$5.25 = 5 \frac{1}{4} = \pm [(1 \cdot 2^2 + 0 \cdot 2^1 + 1) + (0 \cdot 2^{-1} + 1 \cdot 2^{-2})] = \pm (101.01)_2$$

binary expansion

binary

representation

$$(11 \dots 1)_2 = 2^m - 1 \text{ เป็น the largest } m\text{-bit integer}$$

m ตัว

$$\text{และ } (0.11 \dots 1)_2 = 1 - 2^{-k} \text{ เป็น the largest } k\text{-bit fraction}$$

k ตัว

ตัวอย่าง

$(11111)_2 = 2^5 - 1 = 31$ เป็น the largest 5-bit integer
และ $(0.111)_2 = 1 - 2^{-3} = 7/8$ เป็น the largest 3-bit fraction

1.3B การเก็บเลขจำนวนเต็ม (Integers)

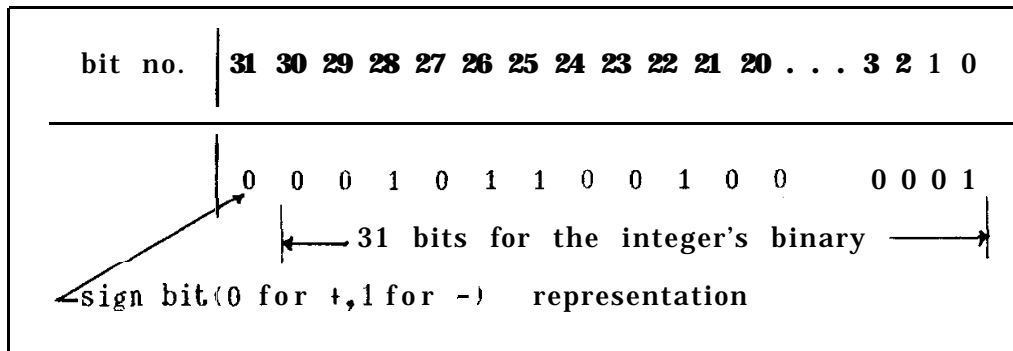
1 word = 32 bits

เลขจำนวนเต็ม X ที่จะเก็บได้ในที่ 32 bits คือ

$$|X| \leq (2^{31} - 1) = 2,147,483,647$$

การพยายามที่จะเก็บเลขจำนวนเต็มที่มีขนาดใหญ่กว่า $(2^{31} - 1)$ จะทำให้เกิด

integer overflow



รูปแสดงการเก็บเลขจำนวนเต็มในขนาด 32 bits

1.3C การเก็บ Floating-Point Representations ของเลขฐานสอง (Binary Numbers)

ทุกๆ nonzero binary number X ถูกแสดงได้ด้วย

$$X = \pm M \cdot 2^c \quad \text{โดยที่ } M = +(0.1b_{-2}b_{-3}\dots b_{-k})_2$$

และ $c =$ เลขจำนวนเต็มตัวหนึ่ง

เรียกว่า normalized floating-point (binary) representation ของ X

ตัวอย่าง

ตาราง 1.3-1 การเขียนในรูปแบบ Normalized Floating-Point

X	Binary Representation	Normalized Binary Floating-Point Representation
-3.5	$(1\ 1.1)_2$	$-(0.1\ 11)_2 \cdot 2^2$, that is, $M = (0.1\ 11)_2$ and $c = 2$
$\frac{7}{8}$	$(0.111)_2$	$(0.1\ 11)_2 \cdot 2^0$, that is, $M = (0.111)_2$ again, but $c = 0$
$\frac{1}{16}$	$(0.0001)_2$	$(0.1)_2 \cdot 2^{-3}$, that is, $M = (0.1)_2$, and $c = -3$

คำว่า Normalized นั้นมาจากการที่ k-bit fraction M (ซึ่งเรียกว่า mantissa) มี $b_{-1} = 1$ ดังนั้น $M \geq 1/2$

ส่วนคำว่า floating point นั้นอ้างถึงการที่เลขจำนวนเต็ม c (ซึ่งเรียกว่า characteristic หรือ exponent ของ X) ถูกปรับค่าเพื่อทำให้ binary point นั้นอยู่ทางซ้ายของ 1 (คือ b_{-1}) ใน binary representation ของ X (นั่นคือเพื่อทำให้ $b_{-1} = 1$ นั่นเอง)

เครื่องคำนวณแบบดิจิทัลทุกเครื่องจะเก็บ mantissa ในรูปที่ normalized แล้ว ทั้งนี้เพื่อให้แน่ใจว่าได้เก็บ bit ของ M มากที่สุดเท่าที่จะเป็นไปได้

ตัวอย่าง

$1/5 = 0.2 = (0.0011001100110011\dots)_2$
 นั้นอาจถูกเก็บใน 6-bit mantissa device ดังนี้

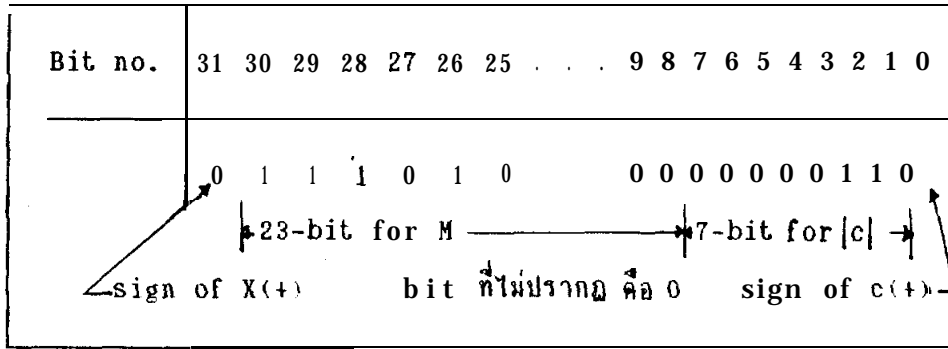
Not normalized: $(0.001100)_2 \cdot 2^0 = 1/8 + 1/16 = 3/16$

$$(\text{error} = 1/5 - 3/16 = 1/80)$$

Normalized: $(0.110011)_2 \cdot 2^{-2} = (1/2 + 1/4 + 1/32 + 1/64) \cdot 1/4 = 51/256$

$$(\text{error} = 1/5 - 51/256 = 1/1280)$$

Normalized representation มีความแม่นยำมากกว่า Not normalized representation ถึง 16 เท่า ทั้งนี้เพราะไม่เสีย 2 bits เพื่อเก็บ (nonsignificant) leading zeros ของ M



รูปแสดงการเก็บ a floating point binary representation ใน a 32-bit word

จากรูป $c = +(0000011)_2 = +3$

$M = (0.11101)_2$ และ sign bit ของ X คือ 0 (สำหรับ +)

ดังนั้น $x = +(0.11101)_2 \cdot 2^3 = +(111.01)_2 = 7.14$

1.3D 23-bit Mantissa Device สามารถเก็บเลขจำนวนขนาดใดได้บ้าง ?

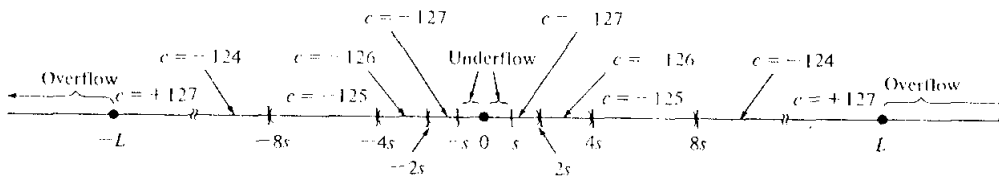
เราจะเรียก nonzero binary number ว่า X representable ถ้าเลขฐานสองจำนวนนั้นสามารถถูกเก็บได้ตามรูปบน นั่นคือถ้า

$X = \pm M \cdot 2^c$ โดยที่ M คือ normalized 23-bit fraction

$$\text{และ } |c| \leq 127 \quad \dots (7)$$

$1/2 = (0.10 \dots 0)_2 \leq M \leq (0.11 \dots 1)_2 = 1 - 2^{-23}$ (ทั้งนี้ b_{-1} ต้องเป็น 1)

ดังนั้น $(1/2)2^c \leq M \cdot 2^c = |X| \leq 1 \cdot 2^c$ (ทั้งนี้เพราะ $M < 1$)



รูป 1.3-1 ช่วงสำหรับ c ค่าต่างๆ กัน

และเนื่องจากว่า $-127 < c < 127$

และ $1/2 < M \leq 1-2^{-23}$

ดังนั้น $s = (1/2)2^{-127} = 2^{-128} = 2.938736... \times 10^{-39}$ คือ smallest representable $X > 0$... (8)

และ $L = (1-2^{-23})2^{127} = 1.70141... \times 10^{28}$ คือ largest representable X ... (9)

ถ้าพยายามเก็บ x ซึ่ง $x > L$ แล้ว จะทำให้เกิด floating-point overflow

ในทำนองเดียวกัน ถ้าพยายามเก็บ x ซึ่ง $0 < x < s$ จะทำให้เกิด floating-point underflow

1.3E ความคลาดเคลื่อนแบบฝังติด (Inherent Error)

เลขจำนวนจริง x ส่วนมากจะไม่สามารถถูกเก็บอย่าง exactly ใน k -bit mantissa device เครื่องจะเก็บค่าประมาณ คือ $X = \pm M \cdot 2^c$ ซึ่งมี error ใน the least significant bit ของ M (นั่นคือ b_{-k}) เราเรียก error นี้ว่า ความคลาดเคลื่อนแบบฝังติดในการเก็บค่า x ซึ่งจะมีค่าน้อยกว่า 2^{-k}

ถ้า $k = 23$ แล้ว M จะมีความคลาดเคลื่อนน้อยกว่า $2^{-23} \approx 1.2 \times 10^{-7}$ ดังนั้น 23-bit mantissa device จะเก็บเลขจำนวนจริง x ได้แม่นยำถึง ประมาณ 7 ตัวเลขนัยสำคัญ (7 significant digits)

ตัวอย่าง

$$x = 0.1 = 0.8 \times 2^{-3}, \quad x = M \cdot 2^{-3}$$

$$M = \overbrace{(0.1100110011001100110)}^{23 \text{ bits}}_2 = 6,710,886/8,388,608 = 0.7999999523...$$

ดังนั้น $x = 0.1$ จะเก็บอยู่ในค่าซึ่งเล็กกว่าค่าจริง

$$x = (0.7999999523\dots) \cdot (1/8) = 0.09999999404\dots$$

*To minimized inherent roundoff, variables x that are input as decimals should be entered rounded to as many significant digits as can possibly be stored on the device used.

จากตัวอย่างการวิงโปรแกรมในหัวข้อ 1.2A นั้น สำหรับเครื่องที่เก็บได้ 8 significant digits เราจะ enter ค่า $\pi = 3.1415927$

1.4 ความคลาดเคลื่อนในการคำนวณแบบ Fixed-Precision

1.4A Decimal-Place UP:: Significant-Digit Accuracy

“leading significant digit” ของเลขทศนิยมคือ “leftmost nonzero digit”

$$\text{ดังนั้น} \quad -0.\overbrace{00666667}^{8d} \text{ approximates } -1/150$$

6s

$$\text{นั่นคือ} \quad -1/150 = -0.00666667(6s)$$

$$\text{หรือ} \quad -1/150 = -0.00666667(8d)$$

$$4 \frac{1}{9} = 4.011111\dots = 4.01(2d \text{ หรือ } 3s)$$

$$4 \frac{1}{9} = 4.011(3d \text{ หรือ } 4s)$$

หรืออาจเขียน

$$-1/150 \doteq -0.00666667$$

$$4 \frac{1}{9} \doteq 4.01 \text{ และ } 4 \frac{1}{9} \doteq 4.011$$

Every digital device has a fixed precision that can be described in terms of the number of significant digits that it can accurately store in one word.

ตัวอย่าง สำหรับ 23-bit mantissa device ที่สามารถเก็บเลขจำนวนได้ด้วย
ความแม่นยำประมาณ 7 ตัวเลขนี้สำคัญ นั่นคือ 7s device

real numbers: x, y, z,... (ใช้ lower case)

ค่า approximate: X, Y, Z,... (ใช้ upper case)

ดังนั้นใน 7s device: $x = X(7s)$

1.4B การคำนวณแบบ Fixed-Precision

พิจารณาเครื่องคิดเลขชนิด 4s, ฐาน 10

$$X = \pm M \cdot 10^c, \quad M = +0.d_1d_2d_3d_4 \quad (d_1 \neq 0 \text{ if } X \neq 0)$$

c เลขจำนวนเต็มซึ่งมีค่าระหว่าง -9 ถึง 9

(Normalized floating-point(decimal) representation)

ดังนั้น $x = -1/150$ จะถูกเก็บเป็น $X = -0.6667 \times 10^{-2}$ ($M = 0.6667, c = -2$)

$y = 4 \frac{1}{9}$ จะถูกเก็บเป็น $Y = 0.4011 \times 10^1$ ($M = 0.4011, c = +1$)

ค่าบวกที่เล็กที่สุดและใหญ่ที่สุดคือ $s = 0.1000 \times 10^{-9} = 10^{-10}$

$$\text{และ } L = 0.9993 \times 10^0 = 999,900,000$$

ให้ "o" แทน t, -, *, /

ในการคำนวณหา $x \circ y$ จะหา $x \hat{o} Y$ แทน โดยที่

(**) $x \hat{o} Y$ หมายความว่า: $W1 \times 0 Y$ อย่างน้อย 5s แล้วปิดเศษให้

เหลือ 4s

เราเรียกวิธีนี้ว่า "4s arithmetic" ซึ่งอาจเป็นผลให้เกิดความคลาดเคลื่อน

หลายๆชนิด

1.4C ความคลาดเคลื่อนในการคำนวณแบบ Fixed-Precision

ban 5 จำนวนคือ

$x=8846.4$, $y=0.0012495$, $z=0.40366$, $u=0.4037681$, $n=50$ และจะถูกเก็บ
ใน 4s device ดังนี้

$$X=0.8846 \times 10^4, Y=0.1250 \times 10^{-2}, Z=0.4037 \times 10^0, W=0.4038 \times 10^0$$

$$\text{และ } U=0.5000 \times 10^2$$

การคำนวณตาม (***) ในหัวข้อ 1.4B จะเกิด errors ต่างๆ เช่น

$$x+y=8846.4012495 \neq x, \quad \text{แต่ } X \uparrow Y = 8846 = X \quad (1)$$

$$z+w=0.8074281 \neq 0.8074, \quad \text{แต่ } Z \uparrow W = Z+W = 0.8075 \quad (2)$$

$$z-y=0.4024105 \neq 0.4024, \quad \text{แต่ } z \hat{-} Y = 0.4025 \quad (3)$$

$$x^2=78,258,793(8s) \neq 0.7826 \cdot 10^8, \quad \text{แต่ } X \hat{*} X = 0.7825 \cdot 10^8 \quad (4)$$

$$u/y=40016.0(6s) \neq 0.4002 \cdot 10^5, \quad \text{แต่ } u \hat{/} Y = 0.4000 \cdot 10^5 \quad (5)$$

$$z-w=-0.0001081 \neq -0.1081 \cdot 10^{-3}, \quad \text{แต่ } Z \hat{-} W = -0.0001 = -0.1000 \cdot 10^{-3} \quad (6)$$

(1)-(5) แสดงความคลาดเคลื่อนแพร่กระจายเนื่องจากการปัดเศษ
(propagated roundoff error)

(1) แสดงการบวกเลขที่มีขนาดใหญ่หลายๆ กับเลขที่มีขนาดเล็กมากๆ ซึ่งผลบวกก็คือ
เลขจำนวนที่มีขนาดใหญ่ของมันเอง

(2)-(5) แสดงว่า ถึงแม้ว่าเราจะปัดเศษ X และ Y ให้มีตัวเลขนัยสำคัญ s ตัว
 $X \hat{o} Y$ ก็อาจเกิดความคลาดเคลื่อนในตัวเลขนัยสำคัญตัวที่ s ได้ ดังนั้น
digital device จึงมักเก็บตัวเลขไว้มากกว่าที่จะแสดงได้

(4)-(5) แสดงว่าความคลาดเคลื่อนจะถูกขยายขนาดขึ้นโดยการคูณด้วยจำนวนที่มี
ขนาดใหญ่ หรือโดยการหารด้วยจำนวนที่มีขนาดเล็ก

$$(X \hat{*} X) - x^2 = -8793 \quad \text{แสดงว่าเกิดความไม่แม่นยำเล็กน้อยใน}$$

$$(U \hat{/} Y) - (u/y) = -16 \quad \text{least significant digit แต่จริงๆ แล้ว ค่าความแตกต่างค่อนข้างใหญ่}$$

เราเรียกความคลาดเคลื่อนนี้ว่าการขยายตัวของความคลาดเคลื่อน

(error magnification) ซึ่งอาจจะเสียหายมากในบางสถานการณ์

- (6) แสดงความคลาดเคลื่อนที่เกิดเมื่อทำการลบเลขสองจำนวนซึ่งมีค่าใกล้เคียงกัน
ใน (6) เมื่อเอา W ไปลบจาก Z ตัวเลขนัยสำคัญสามตัวแรกถูกตัดทิ้งไป
ตัวเลขนัยสำคัญที่เป็นตัวนำของผลลบ (คือของ $Z-W$) มาจากหลักที่ 4 ของ
 W และ Z เท่านั้น แสดงให้เห็นว่าเกิดความคลาดเคลื่อนมากเนื่องจากการ
ปัดเศษ (roundoff) นั่นคือ loss of significance หรือ
catastrophic cancellation

สรุปเกี่ยวกับความคลาดเคลื่อน

- 1) Underflow/Overflow การปฏิบัติการที่ทำให้ได้เลขจำนวนที่มีขนาดเล็กเกินไป
ไปหรือใหญ่เกินไปที่จะถูกเก็บในเครื่องได้
- 2) ความคลาดเคลื่อนจากการปัดเศษ (Roundoff error) เกิดเมื่อทำการคำนวณใน
fixed-precision digital device ซึ่งเป็นผลมาจาก
 - a) ความคลาดเคลื่อนแบบฝังติด (Inherent roundoff) คือการเก็บเลขฐานสอง
ซึ่งเป็นเพียงค่าประมาณของเลขจำนวนจริง
 - b) ความคลาดเคลื่อนแพร่กระจาย (Propagate roundoff) คือการใช้ fixed
precision arithmetic กับเลขจำนวนซึ่งถูกเก็บไว้
 - b.1) Negligible addition การบวกหรือลบเลขที่มีขนาดต่างกันมาก
 - b.2) การขยายตัวของความคลาดเคลื่อน (Error magnification)
การคูณด้วยเลขที่มีขนาดใหญ่หรือ การหารด้วยเลขที่มีขนาดเล็ก
 - b.3) การขาดนัยสำคัญ (Loss of significance) การลบเลขสองจำนวน
ที่มีขนาดใกล้เคียงกันมาก (nearly equal)

1.4D ฟังก์ชันภายใน (Internal Functions)

บิลท์อินฟังก์ชัน (built-in function) เช่น \sin , \cos , \exp , \ln
นั้นโดยทั่วไปแล้วจะทำให้เกิดความคลาดเคลื่อนแพร่กระจายซึ่งเนื่องมาจากการปัดเศษ
มากกว่าที่เกิดจากการคำนวณโดยใช้ตัวดำเนินการทางเลขคณิตตัวหนึ่งๆ (+, -, *, /)

ตัวอย่าง การคำนวณ x^z นั้นจะคำนวณโดยการใช้สูตร

$$x^z = e^{z \ln x} \quad \text{ทั้งนี้โดยที่ฐาน } x \text{ ต้องมีค่าเป็นบวก นั่นคือ}$$

การคำนวณตามคำสั่งในภาษาต่างๆ เช่น

FORTRAN : $X^{**}Z$

BASIC : $X \uparrow Z$

เครื่องคิดเลข : x^z

ถ้า Z เป็นเลขจำนวนเต็มบวก (positive integer) แล้ว $x^z = e^{z \ln x}$ จะแม่นยำน้อยกว่า $X * X * X * \dots * X$ (Z ครั้ง) สำหรับ Z ใดๆ ($2 \leq Z \leq 4$) ความคลาดเคลื่อนแพร่กระจายซึ่งเนื่องมาจากการปัดเศษจะลดลงได้โดยเขียน $X * X * X * \dots * X$ แทน x^z ทั้งนี้ก็เนื่องจากว่า เราทราบว่าคอมพิวเตอร์ทำการคำนวณตามนี้โดยอัตโนมัติอยู่แล้ว

1.5 การหลีกเลี่ยงความคลาดเคลื่อนจากการปัดเศษ (Roundoff Error)

กฎต่างๆไปในการหลีกเลี่ยงการเกิดความคลาดเคลื่อนหรืออย่างน้อยก็ทำให้เกิดน้อยที่สุด

1.5A กฎสำหรับการลดความคลาดเคลื่อนจากการปัดเศษ

กฎที่ 1 เพื่อทำให้ความคลาดเคลื่อนแบบฝังคติน้อยที่สุด จงใช้เลขจำนวนจริง โดยใช้ตัวเลขนัยสำคัญจำนวนเท่ากับที่เครื่องเก็บได้... (1.3E)

กฎที่ 2 เพื่อทำให้โอกาสของการเกิด overflow และ underflow น้อยที่สุด ให้พยายามเก็บค่าที่เกิดจากการคำนวณในทุกขั้นของการคำนวณให้ใกล้ ± 1 (นั่นคือพยายามให้ characteristic มีค่าใกล้ศูนย์)

ตัวอย่าง ถ้าทราบขนาดของ x, y, z แล้ว

นิพจน์ xy/z ควรจะถุกคำนวณดังนี้

$(xy)/z$ ถ้าระหว่าง $|x|$ หรือ $|y|$ มีตัวใดตัวหนึ่งใหญ่และอีกตัวเล็ก

$x(y/z)$ ถ้า $|y|$ และ $|z|$ ทั้งสองตัวเล็ก หรือทั้งสองตัวใหญ่

$(x/z)y$ ถ้า $|x|$ และ $|y|$ ทั้งสองตัวเล็ก หรือทั้งสองตัวใหญ่

ในเครื่องคอมพิวเตอร์ 32-bit word ส่วนมาก overflow (หรือ underflow) จะเกิดขึ้นเมื่อคำนวณ e^x สำหรับ $x > 40$ ($x < -40$) จึงทำให้จำเป็นและสำคัญมากที่จะต้อง ให้ argument ของ exponential function นั้นใกล้ศูนย์มากเท่าที่จะทำได้

ตัวอย่าง ถ้าทราบว่า $x \gg 0, z \gg 0$ และ $n > 1$

ในการคำนวณหรือเขียนคำสั่งในโปรแกรมควรจะทำดังนี้

หาค่าของ e^x/e^z ด้วยคำสั่งหา e^{x-z}

และหาค่าของ z^n/e^{nx} ด้วยคำสั่งหา $(z/e^x)^n$

กฎที่ 3 เพื่อช่วยทำให้ความคลาดเคลื่อนแพร่กระจายน้อยลง ให้ใช้นิพจน์สำหรับการคำนวณที่ต้องใช้การคำนวณทางเลขคณิตน้อยครั้งที่สุด

ตัวอย่าง เช่น $B = \sin^2 x / [x^2(1+\cos x)]$ (*)

ถ้าคำนวณโดยใช้ $(\sin x / x)^2 / (1+\cos x)$ จะลดการคูณลงได้หนึ่งครั้งจากการคำนวณโดยตรงตามสูตร (*)

กฎข้อ 3 นี้จะใช้กับการคำนวณหาค่า polynomial ต่อไปนี้

Nested Multiplication Rule:

Polynomial:

Exponential form: $p_{n \times p}(x) = a_1 x^5 + a_2 x^4 + a_3 x^3 + a_4 x^2 + a_5 x + a_6$

Nested form: $p_{n \times p}(x) = (((((a_1 x + a_2) x + a_3) x + a_4) x + a_5) x + a_6$

$p_{exp}(x)$ นั้นการคำนวณต้องใช้ x^2 subroutine ซึ่งมีความแม่นยำน้อยกว่า (less accurate) และถ้า x เป็นลบจะคำนวณไม่ได้

ตารางต่อไปนี้เป็นเปรียบเทียบความแม่นยำของ $p_{exp}(x)$ และ $p_{nest}(x)$
(ใช้ 23-bit mantissa device)

4th-degree polynomial

$$p_{exp}(x) = x^4 - 9.5x^3 + 28.49x^2 - 28.417x + 2.5662$$

$$p_{nest}(x) = ((x-9.5)x + 28.49)x - 28.417)x + 2.5662$$

$$p_{exp}(x) = p_{nest}(x) = (x-0.1)(x-2.1)(x-2.6)(x-4.7)$$

ตาราง 1.5-1 การเปรียบเทียบความแม่นยำของ $p_{exp}(x)$ และ $p_{nest}(x)$

x	Calculated $p_{exp}(x)$	Calculated $p_{nest}(x)$	Exact Value
2.4	0.317401'	0.317399	0.3174
2.5	0.211199	0.211200*	0.2112
2.7	-0.312001	-0.312000*	-0.3120
2.8	-0.718201	-0.718200*	-0.7182
2.9	-1.209599*	-1.209601	-1.2096
3.1	-2.399997	-2.400001'	-2.4000
3.2	-3.069001	-3.068999"	-3.0690
3.3	-3.763199	-3.763200*	-3.7632
3.5	-5.140801	-5.140800*	-5.1408
3.6	-5.775004	-5.775001'	-5.7750
3.7	-6.335999	-6.336001*	-6.3360
3.8	-6.793202	-6.793201*	-6.7932
3.9	-7.113601*	-7.113602	-7.1136
4.3	-6.283200*	-6.283197	-6.2832
4.4	-5.340597	-5.340602*	-5.3406
4.6	-2.250000*	-2.249999	-2.2500
4.7	0.000010	0.000001*	0.0000
4.8	2.791805	2.791802*	2.7918
4.9	6.182398	6.182401*	6.1824

* แสดงค่าที่แม่นยำกว่า ในตารางแสดงค่า x ซึ่งมีค่าระหว่าง 0.0(0.1)5.0

$$\text{ซึ่ง } |p_{nest}(x) - p_{exp}(x)| > 10^{-6}$$

รากของสมการ polynomial คือ 0.1, 2.1, 2.6 และ 4.7

จากตารางเราจะเห็นว่า

1. $p_{\text{nested}}(x)$ ส่วนมากแม่นยำเท่ากับ $p_{\text{exp}}(x)$
2. ทั้ง $p_{\text{exp}}(x)$ และ $p_{\text{nested}}(x)$ มีความคลาดเคลื่อนใกล้ๆรากตัวที่ใหญ่ที่สุดของสมการ polynomial (นั่นคือ $x = 4.7$)

โดยทั่วไป แม้ว่าจะใช้วิธี nested ก็กับการหา nth-degree polynomial

$$p(x) = a_1x^n + a_2x^{n-1} + a_3x^{n-2} + \dots + a_nx + a_{n+1} \quad (a_1 \neq 0)$$

ซึ่งมีความลำบากที่จะให้ได้ความแม่นยำเมื่อค่าของ x อยู่ใกล้ๆรากจริงที่มีขนาดใหญ่ โดยเฉพาะอย่างยิ่งถ้าสัมประสิทธิ์บางตัวมีขนาดใหญ่ด้วย (เป็นสาเหตุของการขยายตัวของความคลาดเคลื่อน) หรือ slope ใกล้ๆรากมีค่าเข้าใกล้ศูนย์ และยิ่ง n มีค่าสูงขึ้น ๆ สถานการณ์ก็จะยิ่งเลวลง

รูปต่อไปแสดง Synthetic Division Algorithm (Horner's method) เพื่อหาค่า polynomial ในรูป nested

<p>(a) Algorithm: Synthetic Division (Horner's Method)</p> <p>Purpose: To evaluate $PolValue = a_1x^n + a_2x^{n-1} + \dots + a_nx + a_{n+1}$</p> <p>GET $n, a_1, a_2, \dots, a_{n+1}, x$ $PolValue \leftarrow 0$ DO FOR $i = 1$ TO $n + 1$ $PolValue \leftarrow PolValue * x + a_i$</p>																		
<p>(b) Find $p(x) = (((lx - 9.5)x + 28.49)x - 28.417)x + 2.5662$ when $x = -2$</p> <div style="text-align: center;"> <table style="border-collapse: collapse; margin: auto;"> <tr> <td style="border-right: 1px solid black; padding: 5px;">$-2 \mid$</td> <td style="padding: 5px;">1</td> <td style="padding: 5px;">-9.5</td> <td style="padding: 5px;">28.49</td> <td style="padding: 5px;">-28.417</td> <td style="padding: 5px;">2.5662</td> </tr> <tr> <td style="border-right: 1px solid black; padding: 5px; text-align: center;">⊕</td> <td style="padding: 5px;"></td> <td style="padding: 5px;">-2.0</td> <td style="padding: 5px;">23.0</td> <td style="padding: 5px;">-102.98</td> <td style="padding: 5px;">262.794</td> </tr> <tr style="border-top: 1px solid black;"> <td style="border-right: 1px solid black; padding: 5px;"></td> <td style="padding: 5px;">1</td> <td style="padding: 5px;">-11.5</td> <td style="padding: 5px;">51.49</td> <td style="padding: 5px;">-131.397</td> <td style="padding: 5px;">265.3602 = $p(-2)$</td> </tr> </table> </div> <p>Note 1: Diagonal arrows indicate multiplication by $x = -2$.</p> <p>Note 2: $\frac{p(x)}{x - (-2)} = x^3 - 11.5x^2 + 51.49x - 131.397 + \frac{265.3602}{x - (-2)}$</p>	$-2 \mid$	1	-9.5	28.49	-28.417	2.5662	⊕		-2.0	23.0	-102.98	262.794		1	-11.5	51.49	-131.397	265.3602 = $p(-2)$
$-2 \mid$	1	-9.5	28.49	-28.417	2.5662													
⊕		-2.0	23.0	-102.98	262.794													
	1	-11.5	51.49	-131.397	265.3602 = $p(-2)$													

รูป 1.5-1 (a) SYNTHETIC Division Algorithm สำหรับการทำ $p(x)$;
 (b) $p(-2)$; (c) POLVAL


```

(c) 00100.      FUNCTION POLVAL(X, NP1, A)
00200          DIMENSION A(NP1)
00300          C ----- C
00400          C POLVAL = A(1)*X**N + . . . + A(N)*X + A(NP1) C
00500          C WHERE N = DEGREE OF POLYNOMIAL C
00600          C NP1 = NUMBER OF COEFFICIENTS = N+1 C
00700          C ----- C
00800          POLVAL = 0.0
00900          DO 1 I=1, NP1
01000          1 POLVAL = POLVAL*X + A(I)
01100          RETURN
01200          END

```

กฎที่ 4 ออกแบบโปรแกรมเพื่อหลีกเลี่ยงการขาดนัยสำคัญ (loss of significance) ที่อาจเกิดขึ้นเมื่อทำการบวกหรือการลบ

ตัวอย่างใน 1.5B-1.5D จะแสดงการใช้กฎข้อ 4 นี้

1.5B การจัดนิพจน์ตรีโกณมิติเสียใหม่ (Trigonometric Rearrangement) เพื่อหลีกเลี่ยงการขาดนัยสำคัญ (Loss of Significance)

ปัญหา

อธิบายการหาค่า $f(x) = (1 - \cos x) / x^2$ อย่างถูกต้องแม่นยำ สำหรับทุกค่าของ $x \neq 0$

Solution

จาก 1.2A การคำนวณค่าของ $A = (1 - \cos x) / x^2$ นั้นค่าจะไม่แม่นยำเมื่อ $x \approx 0$ เหตุผลเพราะ $(\cos x \approx 1)$ เมื่อ $(x \approx 0)$ เศษจะเกิดการขาดนัยสำคัญ และเมื่อหารด้วย x^2 ซึ่งมีค่าน้อยความคลาดเคลื่อนจะมีมากขึ้น ส่วนการหาค่าของ $B = (\sin x / x)^2 [1 / (1 + \cos x)]$ จะไม่แม่นยำเมื่อ $x \approx \pi$ ทั้งนี้เพราะ $(\cos x \approx -1)$ เมื่อ $(x \approx \pi)$ ทำให้ตัวหารขาดนัยสำคัญและค่าของ B ก็จะไม่เล็ก (คำตอบจะดู absurd)

ดังนั้นเมื่อเข้าใจปัญหาที่เกิดขึ้น เราอาจหลีกเลี่ยงสิ่งที่จะเกิดโดยการคำนวณ $f(x)$ ดังนี้

$$f(x) = \begin{cases} B & \text{ถ้า } (1 - \cos x) \approx 0 \quad (|1 - \cos x| \leq 0.01) \\ A & \text{ถ้า } 1 - \cos x > 0.01 \end{cases}$$

1.5C การจัดนิพจน์พีชคณิตเสียใหม่ (Algebraic Rearrangement)

เพื่อหลีกเลี่ยงการขาดนัยสำคัญ

จาก 1.2B เราพบว่า ถ้า $b_1^2 \gg |c|$ และ $b_1^2 - c_1$ จะไม่ต่างจาก b_1^2 มากดังนั้น $r(-) = b_1 - \sqrt{b_1^2 - c_1}$ จะ loss of significance ถ้า $b_1 > 0$
 $r(+) = b_1 + \sqrt{b_1^2 - c_1}$ จะ loss of significance ถ้า $b_1 < 0$

ถ้าหากจะหารากของ $ax^2 + bx + c = 0$ ($a \neq 0$) โดยที่ไม่ให้ขาดนัยสำคัญ

โดยใช้สูตร $root_1 = (\pm) \left(|b_1| + \sqrt{b_1^2 - c_1} \right)$ และ
 $root_2 = c_1 / root_1$

โดยที่ $b_1 = -b/2a$, $c_1 = c/d$ และ (\pm) คือ $(+)$ นอกจาก $b_1 < 0$

ต่อไปนี้เป็นโปรแกรมย่อยซึ่งเขียนในภาษา FORTRAN เพื่อหารากของสมการ

กำลังสอง $ax^2 + bx + c = 0$, $a \neq 0$ ตามวิธีที่กล่าวข้างต้น

```

00100      SUBROUTINE QRROOTS(A, B, C, ROOT1, ROOT2, COMPLEX, IY, PRINT)
00200      LOGICAL PRINT, COMPLEX
00300      C ----- C
00400      C THIS SUBROUTINE FINDS THE TWO ROOTS OF THE QUADRATIC      C
00500      C      AX**2 + BX + C                                          C
00600      C IF PRINT = TRUE, IT PRINTS THEM ON OUTPUT DEVICE IW.      C
00700      C REAL ROOTS (COMPLEX=FALSE) ARE RETURNED AS ROOT1 AND ROOT2, C
00800      C AND COMPLEX ROOTS (COMPLEX=TRUE) AS ROOT1 +OR- I*ROOT2.   C
00900      C ***** VERSION . 5/1/81 ***** C
01000      B1 = -0.5*B/A
01100      CI = CIA
01200      DISCR = B1*B1 - CI
01300      IF (DISCR .LT. 0.) GOTO 10
01400      C
01500      C REAL ROOTS: ROOT1 AND ROOT2
01600      COMPLEX = .FALSE.
01700      ROOT1 = ABS(B1) + SQRT(DISCR)
01800      IF (B1 .LT. 0.) ROOT1 = -ROOT1
01900      ROOT2 = 0.0
02000      IF (ROOT1 .NE. 0.) ROOT2 = C1/ROOT1
02100      C
02200      IF (PRINT) WRITE(IW,1) ROOT1, ROOT2
02300      1 FORMAT(' REAL ROOTS: ',E14.7,' AND ',E14.7)
02400      RETURN
02500      C
02600      C COMPLEX CONJUGATE ROOTS: ROOT1 +OR- I*ROOT2
02700      10 COMPLEX = .TRUE.
02800      ROOT1 = B1
02900      ROOT2 = SQRT(-DISCR)
03000      C
03100      IF (PRINT) WRITE(IW,2) ROOT1, ROOT2
03200      2 FORMAT(' COMPLEX ROOTS:',E15.7,' +OR- I*(',E14.7,')')
03300      RETURN
03400      C
03500      END

```

รูป 1.5-2 ซึ่บรูกั้ในภาษาฟอร์แทรนชื่อ QRROOTS สำหรับหารากของสมการกำลังสอง

โปรแกรมหลักเพื่อหารากจริงของสมการ $x^2 \pm 100.01 x + 1 = 0$ และ
 รากเชิงซ้อน (complex roots) ของสมการ $x^2 - 2x + 3 = 0$ การคำนวณแม่นยำ
 อย่างน้อยที่สุด 7 ตัวเลขนัยสำคัญ (7s)

```

00100 C ***** TESTQR •
00200 C INTERACTIVE PROGRAM TO TEST SUBROUTINE QROOTS (IW = 5)
00300 LOGICAL COMPLEX
00400 C
00500 WRITE (5,1)
00600 1 FORMAT('OTO FIND ROOTS OF AX**2+BX+C, INPUT A, B, C')
00700 READ (5,*) A, B, C
00800 C
00900 CALL QROOTS (A, B, C, RI, R2, COMPLEX, 5, .TRUE.)
01000 C
01100 STOP
01200 END
    
```

TO FIND ROOTS OF AX**2+BX+C, INPUT A, B, C
1, -100.01, 1

REAL ROOTS: 0.1000000E+03 AND 0.1000000E-01

TO FIND ROOTS OF AX**2+BX+C, INPUT A, B, C
1, +100.01, 1

REAL ROOTS: -0.1000000E+03 AND -0.1000000E-01

TO FIND ROOTS OF AX**2+BX+C, INPUT A, B, C
1, -2, 3

COMPLEX ROOTS: 0.1000000E+01 +0R~ I*(0.1414214E+01)

รูป 1.5-3 โปรแกรมซึ่งเรียกใช้ QROOTS

1.5D การใช้โปรแกรมของแมคคอลลินเพื่อหาลำดับเชิงการขาดนัยสำคัญ

ปัญหา คำนวณ $f(x) = (e^x - 1)/x$ อย่างถูกต้องแม่นยำสำหรับทุก $x \neq 0$

Solution เพราะว่า $e^x \rightarrow 1$ ขณะที่ $x \rightarrow 0$ เศษของ $f(x)$ จะขาดนัยสำคัญ

และจะขยายตัวมากขึ้น (magnified) โดยตัวส่วนที่มีขนาดเล็ก

เมื่อ $x \approx 0$ ปัญหาต่างไปจาก 1.5B เพราะเราไม่อาจหา equivalent
 expression เพื่อช่วยหาลำดับเชิงปัญหาข้างต้น ในสถานการณ์เช่นนี้ เราจะใช้
 โปรแกรมของแมคคอลลินที่เหมาะสมเข้าช่วย

อนุกรมของแมคคอลลินสำหรับ

e^x คือ $[1 + x + x^2/2! + x^3/3! + R_3(x)]$

$$f(x) = (1/x) \{ [1+x+x^2/2!+x^3/3!+R_3(x)] - 1 \} = 1+x/2!+x^2/3!+R_3(x)/x$$

โดยที่ $R_3(x)$ คือ *แบบจำลองของเศษ* (Lagrange for of the remainder)

$$R_3(x) = [\exp(\lambda)/4!] x^4 \quad \text{โดยที่ } \lambda \text{ อยู่ระหว่าง } 0 \text{ และ } x$$

โดยที่ $x < 0.01$ แล้ว $R_3(x)/x$ (ซึ่งคือความคลาดเคลื่อนจากการตัดปลาย (truncation error) ของการประมาณ $f(x)$ ด้วย $1 + x/2! + x^2/3!$)

จะสอดคล้องกับ

$$\left| \frac{R_3(x)}{x} \right| = \frac{[\exp(\lambda)] |x|^3}{24} < \frac{\exp(0.01)(0.01)^3}{24} < 0.00000005$$

จึงรับประกันว่า $1 + x/2! + x^2/3!$ ประมาณ $[\exp(x)-1]/x$ ถึง 7d ถ้า

$$|x| < 0.01$$

ดังนั้นถ้าเราหาค่า $f(x)$ ดังนี้

$$f(x) = \begin{cases} [\exp(x) - 1]/x & \text{ถ้า } |x| \geq 0.01 \\ 1 + (1/2)x[1 + (1/3)x] & \text{ถ้า } |x| < 0.01 \end{cases}$$

แล้ว $f(x)$ จะถูกคำนวณด้วยความแม่นยำของเครื่องเมื่อ $|x| \geq 0.01$ และด้วยความแม่นยำถึง 7d เมื่อ $|x| < 0.01$

1.5E Extended Precision

การขยายที่ที่ใช้เก็บเลขจำนวนจาก 1 เวิร์ด (word) เป็น 2 เวิร์ด นั้นทำให้สามารถเก็บตัวเลขนัยสำคัญได้มากกว่าเดิม ทั้งนี้เพราะเวิร์ดที่ถูกเพิ่มเข้าไบนั้นจะใช้ขยายส่วนที่เก็บ mantissa จึงทำให้เก็บตัวเลขนัยสำคัญได้อย่างน้อยสองเท่าของขนาดปกติ

ตัวอย่าง FORTRAN: กำหนดที่ขนาด 2 เวิร์ด ด้วยคำสั่ง

DOUBLE PRECISION variable list

การคำนวณโดยความเที่ยงตรงสองชั้น (double precision) ใช้เวลาเป็นสองเท่าของการคำนวณโดยความเที่ยงตรงชั้นเดียว (single precision) ดังนั้นจึงควรใช้กฎที่ 1-4 ก่อนเพื่อทำให้การใช้การคำนวณโดยความเที่ยงตรงชั้นเดียวมีประสิทธิภาพ นอกจากนี้ความต้องการความแม่นยำสูง และไม่สามารถใช้กฎทั้ง 4 ได้

1.6 ความคลาดเคลื่อน ความแม่นยำ และ การทดสอบเพื่อตรวจความใกล้เคียง

ในหัวข้อนี้จะศึกษาเทอมต่างๆ ที่ใช้ และหามาตรการเพื่อทำให้ได้ความแม่นยำเท่าที่ต้องการ

1.6A ความคลาดเคลื่อนสัมบูรณ์ (Absolute Error) ความคลาดเคลื่อนสัมพัทธ์ (Relative Error) และ ความคลาดเคลื่อนสัมพัทธ์คิดเป็นร้อยละ (Percent Error)

X เป็นค่าประมาณของค่าจริง x

ความแตกต่าง $\epsilon_x = x - X = (\text{exact value}) - (\text{approximate value})$

คือ ความคลาดเคลื่อนสัมบูรณ์ (absolute error หรือ error) (บางเล่มอาจใช้

$\epsilon_x = X - x$) ดังนั้น $x = X + \epsilon_x$

ถ้า $x \neq 0$ แล้วอัตราส่วน

$r_x = \epsilon_x / x = (x - X) / x$ คือ ความคลาดเคลื่อนสัมพัทธ์ (relative error)

สำหรับการประมาณ x ด้วย X

ในสถานการณ์ที่ทราบค่า r_x แล้วความคลาดเคลื่อนสัมบูรณ์อาจหาได้โดย

$$\epsilon_x = X r_x$$

ความคลาดเคลื่อนสัมพัทธ์ r_x บอกว่าส่วนใดของ x ยังผิดพลาดอยู่

ตัวอย่าง

$$\begin{aligned} |r_x| &= 0.02 \\ |\epsilon_x| &< 0.02 |x| \end{aligned}$$

หรือกล่าวได้ว่า X ประมาณค่า x ด้วยความคลาดเคลื่อนไม่เกิน 2 %

โดยทั่วไป $100 r_x = 100 \epsilon_x / x$ คือ ความคลาดเคลื่อนสัมพัทธ์คิดเป็นร้อยละ

(percent error) ของการประมาณ x ด้วย X

ตัวอย่าง ถ้า $x = 1/30$ ถูกประมาณค่าด้วย $x = 0.033(2s)$ ดังนั้น

$$\epsilon_x = 1/30 - 0.033 = 1/30 - 99/3000 = 1/3000 \text{ และ}$$

$$r_x = (1/3000)/(1/30) = 1/100 = .01$$

ดังนั้น 0.033 ประมาณ 1/30 ด้วยความคลาดเคลื่อนไม่เกิน 1%

1.6B ความเกี่ยวข้องกันของ ϵ และ r ต่อความแม่นยำของ X

$$\text{ถ้า } X = 2.3579 \quad \rightarrow 0.2358(4s)$$

$$\text{หรือ } X = 2.357682 \quad \rightarrow 0.2358(4s)$$

เราต้องหาวิธีบอกคอมพิวเตอร์ให้เลือก X ที่ตรงกับความแม่นยำที่เรา

ต้องการ

การทดสอบความคลาดเคลื่อนสัมบูรณ์ (Absolute Error Test)

d = เลขจำนวนเต็มบวกใด ๆ

$$\left| \epsilon_x \right| = \left| x - X \right| < (1/2) \cdot 10^{-d} \quad \rightarrow \text{จะประกันว่า } X \text{ ประมาณค่า } x \text{ ได้แม่นยำถึงทศนิยมตำแหน่งที่ } d$$

ถ้า $x \neq 0$ แล้วความคลาดเคลื่อนสัมพัทธ์ $r_x = \epsilon_x / x$ สามารถถูกใช้เพื่อให้ได้ความแม่นยำเป็นจำนวนตัวเลขที่สำคัญที่ต้องการ

การทดสอบความคลาดเคลื่อนสัมพัทธ์ (Relative Error Test)

s = เลขจำนวนเต็มบวกใด ๆ

$$\begin{aligned} (1) \quad & \left| r_x \right| < (1/2) \cdot 10^{-s} \quad \rightarrow \text{จะประกันว่า } x \text{ ประมาณค่า } x \text{ ได้แม่นยำถึง } s \text{ ตัว เลขที่สำคัญ} \\ (2) \quad & \text{หรือ } \left| \frac{x - X}{x} \right| < (1/2) \cdot 10^{-s} \end{aligned}$$

พิสูจน์ จาก $x = \pm M \cdot 10^c$

$$\text{โดยที่ } 0.1 \leq M = 0.d_1 d_2 d_3 \dots d_s \dots < 1 \quad \dots (3)$$

ถ้า (1) และ (2) จริงแล้ว จาก (2) และ (3)

$$|\epsilon_x| = |x - \bar{x}| < (1/2) \cdot 10^{-s} |x| = (1/2) \cdot 10^{-s} (M \cdot 10^C) < ((1/2) \cdot 10^{-s}) 10^C$$

แสดงว่า ขนาดของความคลาดเคลื่อนของ $X < 1/2$ ณ ทศนิยมตำแหน่งที่ s ของ M
 ใน (3) (นั่นคือ X มีความแม่นยำถึงจำนวนตัวเลขนัยสำคัญ = s ตัว)

1.7 ขั้นตอนวิธีของการทำซ้ำ (Iterative Algorithms)

ทุกขั้นตอนวิธีของการทำซ้ำสำหรับหาค่าที่ต้องการ (\bar{x}) นั้นจะทำ 2 ขั้นตอน
 ต่อไปนี้

(1) **ขั้นเริ่มต้น**: หา ค่าเดาเริ่มต้นค่าหนึ่ง (initial guess) คือ x_0 ซึ่ง
 จะประมาณค่า \bar{x}

(2) **ขั้นทำซ้ำ**: สำหรับ $k = 0, 1, 2, \dots$ ไปจนกระทั่งการทดสอบเพื่อหยุดการ
 ทำซ้ำเป็นจริง

ใช้ x_k เพื่อหา x_{k+1} ซึ่งเป็นค่าประมาณที่ถูกปรับปรุงแล้ว

เราเรียก x_k ว่า การทำซ้ำครั้งที่ k (k th iterate) และให้

$\epsilon_k = \bar{x} - x_k =$ ความคลาดเคลื่อนของ x_k ดังนั้น ϵ_k คือค่าที่จะบวกกับ x_k
 เพื่อให้ได้ \bar{x} ($\bar{x} = x_k + \epsilon_k$)

ขั้นตอนวิธีจะสำเร็จ (success) ถ้า $x_k \rightarrow \bar{x}$ (นั่นคือ $\epsilon_k \rightarrow 0$) เมื่อ $k \rightarrow \infty$

ขั้นตอนวิธีจะแสดงลักษณะต่อไปนี้

1) Robustness

x_k จะมีค่าลู่เข้า \bar{x} ถึงแม้ว่า x_0 จะไม่ใช่ค่าเดาเริ่มต้นที่ดี

2) Rapid Convergence

เมื่อ x_k เข้าใกล้ \bar{x} , x_{k+1} จะเข้าใกล้มากกว่า

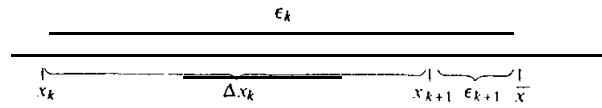
นั่นคือ
$$\left| \epsilon_{k+1} \right| = \left| \bar{x} - x_{k+1} \right| \ll \left| \bar{x} - x_k \right| = \left| \epsilon_k \right|$$

 โดยที่
$$\left| \epsilon_{k+1} \right| \ll \left| \epsilon_k \right|$$

ดังนั้น

$$\Delta x_k = x_{k+1} - x_k \approx \bar{x} - x_k = \epsilon_k \quad (\text{ความคลาดเคลื่อนที่แท้จริงซึ่งไม่ทราบค่า})$$

คือ increment จาก x_k ถึง x_{k+1}

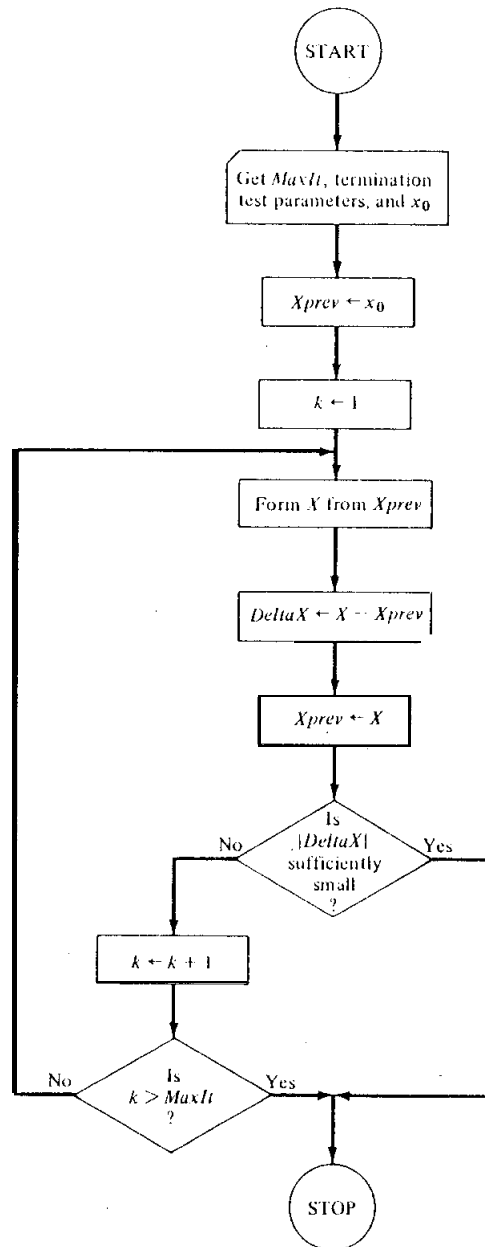


รูป 1.7-1 รูปแสดง Δx_k , ϵ_k , และ ϵ_{k+1}

Δx_k เป็นค่าประมาณที่คำนวณได้ของ ϵ_k

การทดสอบเพื่อหยุดการทำซ้ำ (Termination test) ส่วนมากสำหรับขั้นตอนวิธีของการทำซ้ำซึ่งรวมการหยุดงาน (และใช้ $x_{k+1} \approx \bar{x}$) เมื่อ Δx_k มีค่าเล็กมากพอ

1.7A ผังโปรแกรมสำหรับขั้นตอนวิธีของการทำซ้ำ



รูป 1.7-2 ผังภาพสำหรับ General Iterative Algorithm

General Iterative Algorithm

Purpose: To find \bar{x} to a desired accuracy by repeatedly forming an improved approximation X from a current approximation X_{prev} until either X is sufficiently close to X_{prev} (**termination test**) or $MaxIt$ iterations occur.

(*initialize*)

GET $MaxIt$, termination parameters, x_0 [initial guess]

$X_{prev} \leftarrow x_0$

{ *iterate* }

DO FOR $k = 1$ **TO** $MaxIt$ **UNTIL** **termination test** is satisfied

BEGIN

(form X) $X \leftarrow$ (formula involving X_{prev})

$\Delta X \leftarrow X - X_{prev}$ { ΔX is $\Delta x_{k-1} = x_k - x_{k-1}$ }

[update] $X_{prev} \leftarrow X$ (ready for next iteration)

(**termination test:** $|\Delta X| \approx 0$, as specified by termination parameters)

END

(If **termination test** succeeded, then $\bar{x} \doteq X$ to the desired accuracy. If not, $MaxIt$ iterations failed to yield an X close enough to \bar{x} .)

รูป 1.7-3 รหัสเทียมสำหรับ **General Iterative Algorithm**

1.7B การทดสอบเพื่อหยุดการทำซ้ำ (Termination Tests) สำหรับขั้นตอนวิธีของการทำซ้ำ

จากการประมาณค่า: $\Delta x_k \approx \epsilon_k$ และ ความคลาดเคลื่อนสัมบูรณ์และความคลาดเคลื่อนสัมพัทธ์ (ในข้อ 1.6B) ทำให้เราได้การทดสอบเพื่อหยุดการทำซ้ำ ดังนี้

Absolute Difference Test (For NumDec decimal-place accuracy)

ถ้า $|\Delta X| \leq \text{AbsTol}$ โดยที่ $\text{AbsTol} = \text{Const} * 10^{-\text{NumDec}}$
แล้ว X จะประมาณ \bar{x} ได้แม่นยำถึงทศนิยมตำแหน่งที่ NumDec

Relative Difference Test (For NumSig significant-digit Accuracy)

ถ้า $\frac{\Delta X}{X} \leq \text{RelTol}$ โดยที่ $\text{RelTol} = \text{Const} * 10^{-\text{NumSig}}$
แล้ว X จะประมาณ \bar{x} ได้แม่นยำถึงตัวเลขนัยสำคัญเท่ากับ NumSig ตัว

$$\Delta X = X - X_{\text{prev}}$$

Const เรามักใช้ค่า 1 ถ้าการลู่เข้า (convergence) ช้า จะเลือกค่าเข้าใกล้ 0.5 และถ้าคาดว่าจะมีความคลาดเคลื่อนจากการปัดเศษมาก จะเลือกค่ามากกว่า 1 (คือ 2-9)

ถ้าเราต้องการความแม่นยำถึง 5 ตัวเลขนัยสำคัญ $\bar{x} \approx 3000$ เราควรใช้ NumDec = 1 ใน Absolute Difference Test แต่ถ้าทราบว่า $\bar{x} \approx -0.003$ เราควรใช้ NumDec = 7 ส่วนใน Relative Difference Test เราจะใช้ NumSig = 5 ในทั้ง 2 กรณี นั้นแสดงว่า Relative Difference Test มักจะถูกใช้ในโปรแกรมซึ่งจะหาค่า \bar{x} 's ที่มีขนาดต่างๆกันไป

อย่างไรก็ดี Relative Difference Test นั้นไม่ใช้กับการทดสอบการเข้าใกล้ศูนย์ อันที่จริงแล้ว ถ้า $\bar{x} = 0$ และ ขั้นตอนวิธีลู่เข้าเร็ว (i.e., $x_n \rightarrow 0$ อย่างเร็ว) แล้ว $\left| \frac{\Delta X}{X} \right| \ll \left| \frac{\Delta X}{X_{\text{prev}}} \right|$ ผลที่ตามมาคือ ถ้าไม่มีความคลาดเคลื่อนจากการปัดเศษแล้ว $\left| \frac{\Delta X}{X} \right| \approx \left| \frac{\Delta X}{X_{\text{prev}}} \right|$ และอสมการใน Relative Difference Test จะไม่เป็นจริง

คอมพิวเตอร์อาจพบ ลูปอนันต์ (infinite loop) อย่างที่เราไม่คาดฝัน เราควรระวังกันโดยยึดหลักต่อไปนี้

For **any** iterative algorithm, no matter how sure you are that a termination test will be set, put an upper limit (e.g., MaxIt) on the number of iterations (just in case you are wrong).

1.7C การลู่เข้าแบบเชิงเส้น (Linear Convergence) และ การลู่เข้าแบบอันดับสอง (Quadratic Convergence)

$$(*) \quad (Dx/\text{preceding } Dx) \approx \text{Constant} = C \quad \text{นั่นคือ } \Delta x_k \approx C \Delta x_{k-1}$$

ถ้าเป็นไปตาม (*) นั่นคือเมื่อ $(x_{k+1} - x_k)$ is (approximately) proportional to $(x_k - x_{k-1})$ เราเรียกการลู่เข้าของ x_k to \bar{x} ว่า การลู่เข้าแบบเชิงเส้น (linear convergence)

เหตุการณ์ข้างต้นตรงข้ามกับเมื่อการทำซ้ำ satisfied $(Dx/\text{preceding } Dx) \rightarrow 0$ เมื่อ k เพิ่มขึ้น

$$(**) \quad (Dx/(\text{preceding } Dx))^2 \approx \text{Constant} = C \quad \text{นั่นคือ } \Delta x_k \approx C(\Delta x_{k-1})^2$$

ถ้าเป็นไปตาม (**) นั่นคือเมื่อ $(x_{k+1} - x_k)$ is (approximately) proportional to $(x_k - x_{k-1})^2$ เราเรียกการลู่เข้าของ x_k to \bar{x} ว่า การลู่เข้าแบบอันดับสอง (quadratic convergence)

เราจะใช้ การลู่เข้าแบบอันดับสอง เพื่อเป็นมาตรฐานของเราสำหรับการลู่เข้าอย่างรวดเร็ว (rapid convergence)

และจะใช้ การลู่เข้าแบบเชิงเส้น เพื่อเป็นมาตรฐานของเราสำหรับการลู่เข้าอย่างช้า (slow convergence)