

Accuracy of Numerical Calculations

Misconceptions About Errors

Error is not sin, nor is it always a mistake. Numerical error is merely the difference between what you wish to calculate and what you get. The difference matters only if it is too big. Usually it is negligible; but sometimes error is distressingly big, hard to explain, and harder to correct. This appendix focuses on errors, especially those that might be large—however rare. Here are some examples.

Example 1: A Broken Calculator. Since $(\sqrt{x})^2 = x$ whenever $x \geq 0$, we expect also

$$f(x) = \underbrace{\left(\left(\left(\sqrt{\sqrt{\dots \sqrt{\sqrt{x}}}} \right)^2 \right)^2 \dots \right)^2}_{50 \text{ roots}} \underbrace{\left(\dots \right)^2}_{50 \text{ squares}}$$

should equal x too.

A program of 100 steps can evaluate the expression $f(x)$ for any positive x . When $x = 10$ the HP-15C calculates 1 instead. The error $10 - 1 = 9$ appears enormous considering that only 100 arithmetic operations were performed, each one presumably correct to 10 digits. What the program actually delivers instead of $f(x) = x$ turns out to be

$$f(x) = \begin{cases} 1 & \text{for } x \geq 1 \\ 0 & \text{for } 0 \leq x < 1, \end{cases}$$

which seems very wrong. Should this calculator be repaired?

Example 2: Many Pennies. A corporation retains Susan as a scientific and engineering consultant at a fee of one penny per second for her thoughts, paid every second of every day for a year. Rather than distract her with the sounds of pennies dropping, the corporation proposes to deposit them for her into a bank account in which interest accrues at the rate of $11\frac{1}{4}$ percent per annum compounded every second. At year's end these pennies will accumulate to a sum

$$\text{total} = (\text{payment}) \times \frac{(1 + i/n)^n - 1}{i/n}$$

where payment = \$0.01 = one penny per second,

$i = 0.1125 = 11.25$ percent per annum interest rate,

$n = 60 \times 60 \times 24 \times 365 =$ number of seconds in a year.

Using her HP-15C, Susan reckons that the total will be \$376,877.67. But at year's end the bank account is found to hold \$333,783.35. Is Susan entitled to the \$43,094.32 difference?

In both examples the discrepancies are caused by rounding errors that could have been avoided. This appendix explains how.

The war against error begins with a salvo against wishful thinking, which might confuse what we want with what we get. To avoid confusion, the true and calculated results must be given different names even though their difference may be so small that the distinction seems pedantic.

Example 3: Pi. The constant $\pi = 3.1415926535897932384626433\dots$. Pressing the $\boxed{\pi}$ key on the HP-15C delivers a different value

$$\boxed{\pi} = 3.141592654$$

which agrees with π to 10 significant digits. But $\boxed{\pi} \neq \pi$, so we should not be surprised when, in Radians mode, the calculator doesn't produce $\sin \boxed{\pi} = 0$.

Suppose we wish to calculate x but we get X instead. (This convention is used throughout this appendix.) The *error* is $x - X$. The *absolute error* is $|x - X|$. The *relative error* is usually reckoned $(x - X)/x$ for $x \neq 0$.

Example 4: A Bridge Too Short. The lengths in meters of three sections of a cantilever bridge are designed to be

$$x = 333.76 \quad y = 195.07 \quad z = 333.76 .$$

The measured lengths turn out to be respectively

$$X = 333.69 \quad Y = 195.00 \quad Z = 333.72 .$$

The discrepancy in total length is

$$d = (x + y + z) - (X + Y + Z) = 862.59 - 862.41 = 0.18 .$$

Ed, the engineer, compares the discrepancy d with the total length $(x + y + z)$ and considers the relative discrepancy

$$d/(x + y + z) = 0.0002 = 2 \text{ parts in } 10,000$$

to be tolerably small. But Rhonda, the riveter, considers the absolute discrepancy $|d| = 0.18$ meters (about 7 inches) much too large for her liking; some powerful stretching will be needed to line up the bridge girders before she can rivet them together. Both see the same discrepancy d , but what looks negligible to one person can seem awfully big to another.

Whether large or small, errors must have sources which, if understood, usually permit us to compensate for the errors or to circumvent them altogether. To understand the distortions in the girders of a bridge, we should learn about structural engineering and the theory of elasticity. To understand the errors introduced by the very act of computation, we should learn how our calculating instruments work and what are their limitations. These are details most of us want not to know, especially since a well-designed calculator's rounding errors are always nearly minimal and therefore appear insignificant when they are introduced. But when on rare occasions they conspire to send a computation awry, they must be reclassified as "significant" after all.

Example 1 Explained. Here $f(x) = s(r(x))$, where

$$r(x) = \underbrace{\sqrt{\sqrt{\dots \sqrt{\sqrt{x}}}}}_{50 \text{ roots}} = x^{(1/2^{50})}$$

and

$$s(r) = \underbrace{((\dots ((r)^2)^2 \dots)^2)^2}_{50 \text{ squares}} = r^{(2^{50})}.$$

The exponents are $1/2^{50} = 8.8818 \times 10^{-16}$ and $2^{50} = 1.1259 \times 10^{15}$. Now, x must lie between 10^{-99} and $9.999 \dots \times 10^{99}$ since no positive numbers outside that range can be keyed into the calculator. Since r is an increasing function, $r(x)$ lies between

$$r(10^{-99}) = 0.999999999999997975 \dots$$

and

$$r(10^{100}) = 1.00000000000002045 \dots$$

This suggests that $R(x)$, the calculated value of $r(x)$, would be 1 for all valid calculator arguments x . In fact, because of roundoff,

$$R(x) = \begin{cases} 0.9999999999 & \text{for } 0 < x < 1 \\ 1.0000000000 & \text{for } 1 \leq x \leq 9.999999999 \times 10^{99}. \end{cases}$$

If $0 < x < 1$, then $x \leq 0.9999999999$ in a 10-digit calculator. We would then rightly expect that $\sqrt{x} \leq \sqrt{0.9999999999}$, which is $0.99999999994999999998\dots$, which rounds to 0.9999999999 again. Therefore, if \sqrt{x} is pressed arbitrarily often starting with $x < 1$, the result cannot exceed 0.9999999999 . This explains why we obtain $R(x) = 0.9999999999$ for $0 < x < 1$ above. When $R(x)$ is squared 50 times to produce $F(x) = S(R(x))$, the result is clearly 1 for $x \geq 1$, but why is $F(x) = 0$ for $0 \leq x < 1$? When $x < 1$,

$$s(R(x)) \leq s(0.9999999999) = (1 - 10^{-10})^{2^{50}} \approx 6.14 \times 10^{-48898}.$$

This value is so small that the calculated value $F(x) = S(R(x))$ underflows to 0. So the HP-15C isn't broken; it is doing the best that can be done with 10 significant digits of precision and 2 exponent digits.

We have explained example 1 using no more information about the HP-15C than that it performs each arithmetic operation \sqrt{x} and x^2 fully as accurately as is possible within the limitations of 10 significant digits and 2 exponent digits. The rest of the information we needed was mathematical knowledge about the functions f , r , and s . For instance, the value $r(10^{100})$ above was evaluated as

$$\begin{aligned} r(10^{100}) &= (10^{100})^{(1/2)^{50}} \\ &= \exp(\ln(10^{100})/2^{50}) \\ &= \exp(100(\ln 10)/2^{50}) \\ &= \exp(2.045 \times 10^{-13}) \\ &= 1 + (2.045 \times 10^{-13}) + \frac{1}{2}(2.045 \times 10^{-13})^2 + \dots \end{aligned}$$

by using the series $\exp(z) = 1 + z + \frac{1}{2}z^2 + \frac{1}{6}z^3 + \dots$.

Similarly, the binomial theorem was used for

$$\begin{aligned} \sqrt{0.9999999999} &= (1 - 10^{-10})^{1/2} \\ &= 1 - \frac{1}{2}(10^{-10}) - \frac{1}{8}(10^{-10})^2 - \dots \end{aligned}$$

These mathematical facts lie well beyond the kind of knowledge that might have been considered adequate to cope with a calculation containing only a handful of multiplications and square roots. In this respect, example 1 illustrates an unhappy truism: Errors make computation very much harder to analyze. That is why a well-designed calculator, like the HP-15C, will introduce errors of its own as sparingly as is possible at a tolerable cost. Much more error than that would turn an already difficult task into something hopeless.

Example 1 should lay two common *misconceptions* to rest:

- Rounding errors can overwhelm a computation only if vast numbers of them accumulate.
- A few rounding errors can overwhelm a computation only if accompanied by massive cancellation.

Regarding the first misconception, example 1 would behave in the same perverse way if it suffered only one rounding error, the one that produces $R(x) = 1$ or 0.9999999999 , in error by less than one unit in its last (10th) significant digit.

Regarding the second misconception, cancellation is what happens when two nearly equal numbers are subtracted. For example, calculating

$$c(x) = (1 - \cos x)/x^2$$

in Radians mode for small values of x is hazardous because of cancellation. Using $x = 1.2 \times 10^{-5}$ and rounding results to 10 digits,

$$\cos x = 0.9999999999$$

and

$$1 - \cos x = 0.0000000001$$

with cancellation leaving maybe one significant digit in the numerator. Also

$$x^2 = 1.44 \times 10^{-10}.$$

Then

$$C(x) = 0.6944.$$

This calculated value is wrong because $0 \leq c(x) < 1/2$ for all $x \neq 0$. To avoid numerical cancellation, exploit the trigonometric identity $\cos x = 1 - 2 \sin^2(x/2)$ to cancel the 1 *exactly* and obtain a better formula

$$c(x) = \frac{1}{2} \left(\frac{\sin(x/2)}{x/2} \right)^2.$$

When this latter expression is evaluated (in Radians mode) at $x = 1.2 \times 10^{-5}$, the computed result $C(x) = 0.5$ is correct to 10 significant digits. This example, while explaining the meaning of the word “cancellation,” suggests that it is always a bad thing. That is another misconception to be dispatched later. For the

present, recall that example 1 contains no subtraction, therefore no cancellation, and is still devastated by its rounding error. In this respect example 1 is counterintuitive, a little bit scary. Nowhere in it can we find one or two arithmetic operations to blame for the catastrophe; no small rearrangement will set everything right as happened for $c(x)$. Alas, example 1 is not an isolated example. As computers and calculators grow in power, so do instances of insidious error growth become more common.

To help you recognize error growth and cope with it is the ultimate goal of this appendix. We shall start with the simplest kinds of errors and work our way up gradually to the subtle errors that can afflict the sophisticated computations possible on the HP-15C.

A Hierarchy of Errors

Some errors are easier to explain and to tolerate than others. Therefore, the functions delivered by single keystrokes on the HP-15C have been categorized, for the purposes of easier exposition, according to how difficult their errors are to estimate. The estimates should be regarded as goals set by the calculator's designers rather than as specifications that guarantee some stated level of accuracy. On the other hand, the designers believe they can prove mathematically that their accuracy goals have been achieved, and extensive testing has produced no indication so far that they might be mistaken.

Level 0: No Error

Functions which should map small integers (smaller than 10^{10}) to small integers do so exactly, without error, as you might expect.

Examples:

$$\sqrt{4} = 2 \quad -2^3 = -8 \quad 3^{20} = 3,486,784,401$$

$$\log(10^9) = 9 \quad 6! = 720$$

$$\cos^{-1}(0) = 90 \text{ (in Degrees mode)}$$

$$\text{ABS}(4,684,660 + 4,684,659i) = 6,625,109 \text{ (in Complex mode)}$$

Also exact for real arguments are **ABS**, **FRAC**, **INT**, **RND**, and comparisons (such as $x \leq y$). But the matrix functions **×**, **÷**, **1/x**, **MATRIX**6, and **MATRIX**9 (determinant) are exceptions (refer to page 192).

Level ∞ : Overflow/Underflow

Results which would lie closer to zero than 10^{-99} underflow quietly to zero. Any result that would lie beyond the overflow thresholds $\pm 9.999999999 \times 10^{99}$ is replaced by the nearest threshold, and then flag 9 is set and the display blinks. (Pressing **ON** **ON** or **CF** 9 or **←** will clear flag 9 and stop the blinking.) Most functions that result in more than one component can tolerate overflow/underflow in one component without contaminating the other; examples are **→R**, **→P**, complex arithmetic, and most matrix operations. The exceptions are matrix inversion (**1/x** and **÷**), **MATRIX** 9 (determinant), and **L.R.**.

Level 1: Correctly Rounded, or Nearly So

Operations that deliver “correctly rounded” results whose error cannot exceed $\frac{1}{2}$ unit in their last (10th) significant digit include the real algebraic operations **+**, **-**, **×**, **÷**, **x²**, **√x**, **1/x**, and **%**, the complex and matrix operations **+** and **-**, matrix by scalar operations **×** and **÷** (excluding division by a matrix), and **→H.MS**. These results are the best that 10 significant digits can represent, as are familiar constants **π**, 1 **e^x**, 2 **LN**, 10 **LN**, 1 **→RAD**, and many more. Operations that can suffer a slightly larger error, but still significantly smaller than one unit in the 10th significant digit of the result, include **Δ%**, **→H**, **→RAD**, **→DEG**, **P_{y,x}**, and **C_{y,x}**; **LN**, **LOG**, **10^x**, and **TANH** for real arguments; **→P**, **SIN⁻¹**, **COS⁻¹**, **TAN⁻¹**, **SINH⁻¹**, **COSH⁻¹**, and **TANH⁻¹** for real and complex arguments; **ABS**, **√x**, and **1/x** for complex arguments; matrix norms **MATRIX** 7 and **MATRIX** 8; and finally **SIN**, **COS**, and **TAN** for real arguments in Degrees and Grads modes (but not in Radians mode—refer to Level 2, page 184).

A function that grows to ∞ or decays to 0 exponentially fast as its argument approaches $\pm\infty$ may suffer an error larger than one unit in its 10th significant digit, but only if its magnitude is smaller than 10^{-20} or larger than 10^{20} ; and though the relative error gets worse as the result gets more extreme (small or large), the error stays below three units in the last (10th) significant digit. The reason for this error is explained later. Functions so affected are **e^x**, **y^x**, **x!** (for noninteger x), **SINH**, and **COSH** for real arguments. The worst case known is 3^{201} , which is calculated as $7.968419664 \times 10^{95}$. The last digit 4 should be 6 instead, as is the case for $7.29^{33.5}$, calculated as $7.968419666 \times 10^{28}$.

The foregoing statements about errors can be summarized for all functions in Level 1 in a way that will prove convenient later:

Attempts to calculate a function f in Level 1 produce instead a computed value $F = (1 + \epsilon)f$ whose relative error ϵ , though unknown, is very small:

$$|\epsilon| < \begin{cases} 5 \times 10^{-10} & \text{if } F \text{ is correctly rounded} \\ 1 \times 10^{-9} & \text{for all other functions } F \text{ in Level 1.} \end{cases}$$

This simple characterization of all the functions in Level 1 fails to convey many other important properties they all possess, properties like

- Exact integer values: mentioned in Level 0.
- Sign symmetry: $\sinh(-x) = -\sinh(x)$, $\cosh(-x) = \cosh(x)$, $\ln(1/x) = -\ln(x)$ (if $1/x$ is computed exactly).
- Monotonicity: if $f(x) \geq f(y)$, then computed $F(x) \geq F(y)$.

These additional properties have powerful implications; for instance, $\text{TAN}(20^\circ) = \text{TAN}(200^\circ) = \text{TAN}(2,000^\circ) = \dots = \text{TAN}(2 \times 10^{99}^\circ) = 0.3639702343$ correctly. But the simple characterization conveys most of what is worth knowing, and that can be worth money.

Example 2 Explained. Susan tried to calculate

$$\text{total} = \text{payment} \times \frac{(1 + i/n)^n - 1}{i/n}$$

where

$$\text{payment} = \$0.01,$$

$$i = 0.1125, \text{ and}$$

$$n = 60 \times 60 \times 24 \times 365 = 31,536,000.$$

She calculated \$376,877.67 on her HP-15C, but the bank's total was \$333,783.35, and this latter total agrees with the results calculated on good, modern financial calculators like the HP-12C, HP-37E, HP-38E/38C, and HP-92. Where did Susan's calculation go awry? No severe cancellation, no vast accumulation of errors; just one rounding error that grew insidiously caused the damage:

$$i/n = 0.000000003567351598$$

$$1 + i/n = 1.000000004$$

when rounded to 10 significant digits. There is the rounding error that hurts. Subsequently attempting to calculate $(1 + i/n)^n$, Susan must get instead $(1.000000004)^{31,536,000} = 1.134445516$, which is wrong in its second decimal place.

How can the correct value be calculated? Only by not throwing away so many digits of i/n . Observe that

$$(1 + i/n)^n = e^{n \ln(1 + i/n)},$$

so we might try to calculate the logarithm in some way that does not discard those precious digits. An easy way to do so on the HP-15C does exist.

To calculate $\lambda(x) = \ln(1 + x)$ accurately for all $x > -1$, even if $|x|$ is very small:

1. Calculate $u = 1 + x$ rounded.
2. Then

$$\lambda(x) = \begin{cases} x & \text{if } u = 1 \\ \ln(u) x / (u - 1) & \text{if } u \neq 1. \end{cases}$$

The following program calculates $\lambda(x) = \ln(1 + x)$.

Keystrokes	Display	
g P/R		
f CLEAR PRGM	000-	
f LBL A	001-42,21,11	Assumes x is in X-register.
ENTER	002- 36	
ENTER	003- 36	
EEX	004- 26	Places 1 in X-register.
+	005- 40	Calculates $u = 1 + x$ rounded.
g LN	006- 43 12	Calculates $\ln(u)$ (zero for $u = 1$).
x \leftrightarrow y	007- 34	Restores x to X-register.
g LSTx	008- 43 36	Recalls u .

Keystrokes	Display	
EEX	009-	26 Places 1 in X-register.
g TEST 6	010-43,30, 6	Tests $u \neq 1$.
-	011-	30 Calculates $u - 1$ when $u \neq 1$.
÷	012-	10 Calculates $x/(u - 1)$ or $1/1$.
x	013-	20 Calculates $\lambda(x)$.
g RTN	014- 43 32	
g P/R		

The calculated value of u , correctly rounded by the HP-15C, is $u = (1 + \epsilon)(1 + x)$, where $|\epsilon| < 5 \times 10^{-10}$. If $u = 1$, then

$$|x| = |1/(1 + \epsilon) - 1| \leq 5 \times 10^{-10}$$

too, in which case the Taylor series $\lambda(x) = x(1 - \frac{1}{2}x + \frac{1}{3}x^2 - \dots)$ tells us that the correctly rounded value of $\lambda(x)$ must be just x . Otherwise, we shall calculate $x\lambda(u - 1)/(u - 1)$ fairly accurately instead of $\lambda(x)$. But $\lambda(x)/x = 1 - \frac{1}{2}x + \frac{1}{3}x^2 - \dots$ varies very slowly, so slowly that the absolute error $\lambda(x)/x - \lambda(u - 1)/(u - 1)$ is no worse than the absolute error $x - (u - 1) = -\epsilon(1 + x)$, and if $x \leq 1$, this error is negligible relative to $\lambda(x)/x$. When $x > 1$, then $u - 1$ is so nearly x that the error is negligible again; $\lambda(x)$ is correct to nine significant digits.

As usual in error analyses, the explanation is far longer than the simple procedure being explained and obscures an important fact: the errors in $\ln(u)$ and $u - 1$ were ignored during the explanation because we knew they would be negligible. This knowledge, and hence the simple procedure, is *invalid* on some other calculators and big computers! Machines do exist which calculate $\ln(u)$ and/or $1 - u$ with small *absolute* error, but large *relative* error when u is near 1; on those machines the foregoing calculations must be wrong or much more complicated, often both. (Refer to the discussion under Level 2 for more about this.)

Back to Susan's sum. By using the foregoing simple procedure to calculate $\lambda(i/n) = \ln(1 + i/n) = 3.567351591 \times 10^{-9}$, she obtains a better value:

$$(1 + i/n)^n = e^{n\lambda(i/n)} = 1.119072257$$

from which the correct total follows.

To understand the error in 3^{201} , note that this is calculated as $e^{201 \ln(3)} = e^{220.821\dots}$. To keep the final relative error below one unit in the 10th significant digit, $201 \ln(3)$ would have to be calculated with an absolute error rather smaller than 10^{-10} , which would entail carrying at least 14 significant digits for that intermediate value. The calculator does carry 13 significant digits for certain intermediate calculations of its own, but a 14th digit would cost more than it's worth.

Level 1C: Complex Level 1

Most complex arithmetic functions cannot guarantee 9 or 10 correct significant digits in each of a result's real and imaginary parts separately, although the result will conform to the summary statement about functions in Level 1 provided f , F , and ϵ are interpreted as complex numbers. In other words, every complex function f in Level 1C will produce a calculated complex value $F = (1 + \epsilon)f$ whose small complex relative error ϵ must satisfy $|\epsilon| < 10^{-9}$. The complex functions in Level 1C are $\boxed{\times}$, $\boxed{\div}$, $\boxed{x^2}$, $\boxed{\text{LN}}$, $\boxed{\text{LOG}}$, $\boxed{\text{SIN}^{-1}}$, $\boxed{\text{COS}^{-1}}$, $\boxed{\text{TAN}^{-1}}$, $\boxed{\text{SINH}^{-1}}$, $\boxed{\text{COSH}^{-1}}$, and $\boxed{\text{TANH}^{-1}}$. Therefore, a function like $\lambda(z) = \ln(1 + z)$ can be calculated accurately for all z by the same program as given above and with the same explanation.

To understand why a complex result's real and imaginary parts might not individually be correct to 9 or 10 significant digits, consider $\boxed{\times}$, for example: $(a + ib) \times (c + id) = (ac - bd) + i(ad + bc)$ ideally. Try this with $a = c = 9.999999998$, $b = 9.999999999$, and $d = 9.999999997$; the exact value of the product's real part ($ac - bd$) should then be

$$\begin{aligned} (9.999999998)^2 - (9.999999999)(9.999999997) \\ &= 99.999999980000000004 - 99.999999980000000003 \\ &= 10^{-18} \end{aligned}$$

which requires that at least 20 significant digits be carried during the intermediate calculation. The HP-15C carries 13 significant digits for internal intermediate results, and therefore obtains 0 instead of 10^{-18} for the real part, but this error is negligible compared to the imaginary part 199.9999999 .

Level 2: Correctly Rounded for Possibly Perturbed Input

Trigonometric Functions of Real Radian Angles

Recall example 3, which noted that the calculator's $\boxed{\pi}$ key delivers an approximation to π correct to 10 significant digits but still slightly different from π , so $0 = \sin(\pi) \neq \sin(\boxed{\pi})$ for which the calculator delivers

$$\boxed{\text{SIN}}(\boxed{\pi}) = -4.100000000 \times 10^{-10}.$$

This computed value is not quite the same as the true value

$$\sin(\boxed{\pi}) = -4.10206761537356... \times 10^{-10}.$$

Whether the discrepancy looks small (absolute error less than 2.1×10^{-13}) or relatively large (wrong in the fourth significant digit) for a 10-significant-digit calculator, the discrepancy deserves to be understood because it foreshadows other errors that look, at first sight, much more serious.

Consider

$$10^{14}\pi = 314159265358979.3238462643...$$

with $\sin(10^{14}\pi) = 0$ and

$$10^{14} \times \boxed{\pi} = 314159265400000$$

with $\boxed{\text{SIN}}(10^{14}\boxed{\pi}) = 0.7990550814$, although the true

$$\sin(10^{14}\boxed{\pi}) = -0.78387....$$

The wrong sign is an error too serious to ignore; it seems to suggest a defect in the calculator. To understand the error in trigonometric functions we must pay attention to small differences among π and two approximations to π :

true	$\pi = 3.1415926535897932384626433...$	
key	$\boxed{\pi} = 3.141592654$	(matches π to 10 digits)
internal p	$= 3.141592653590$	(matches π to 13 digits)

Then all is explained by the following formula for the calculated value: $\boxed{\text{SIN}}(x) = \sin(x\pi/p)$ to within ± 0.6 units in its last (10th) significant digit.

More generally, if $\text{trig}(x)$ is any of the functions $\sin(x)$, $\cos(x)$, or $\tan(x)$, evaluated in real Radians mode, the HP-15C produces

$$\boxed{\text{TRIG}}(x) = \text{trig}(x\pi/p)$$

to within ± 0.6 units in its 10th significant digit.

This formula has important practical implications:

- Since $\pi/p = 1 - 2.0676... \times 10^{-13}/p = 0.9999999999999342...$, the value produced by $\boxed{\text{TRIG}}(x)$ differs from $\text{trig}(x)$ by no more than can be attributed to two perturbations: one in the 10th significant digit of the output $\text{trig}(x)$, and one in the 13th significant digit of the input x .

If x has been calculated and rounded to 10 significant digits, the error inherited in its 10th significant digit is probably orders of magnitude bigger than $\boxed{\text{TRIG}}$'s second perturbation in x 's 13th significant digit, so this second perturbation can be ignored unless x is regarded as known or calculated exactly.

- Every trigonometric identity that does not explicitly involve π is satisfied to within roundoff in the 10th significant digit of the calculated values in the identity. For instance,

$$\sin^2(x) + \cos^2(x) = 1, \text{ so } (\boxed{\text{SIN}}(x))^2 + (\boxed{\text{COS}}(x))^2 = 1$$

$$\sin(x)/\cos(x) = \tan(x), \text{ so } \boxed{\text{SIN}}(x)/\boxed{\text{COS}}(x) = \boxed{\text{TAN}}(x)$$

with each calculated result correct to nine significant digits for all x . Note that $\boxed{\text{COS}}(x)$ vanishes for no value of x representable exactly with just 10 significant digits. And if $2x$ can be calculated exactly given x ,

$$\sin(2x) = 2\sin(x)\cos(x), \text{ so } \boxed{\text{SIN}}(2x) = 2\boxed{\text{SIN}}(x)\boxed{\text{COS}}(x)$$

to nine significant digits. Try the last identity for $x = 52174$ radians on the HP-15C:

$$\boxed{\text{SIN}}(2x) = -0.00001100815000,$$

$$2\boxed{\text{SIN}}(x)\boxed{\text{COS}}(x) = -0.00001100815000.$$

Note the close agreement even though for this x , $\sin(2x) = 2\sin(x)\cos(x) = -0.0000110150176...$ disagrees with $\boxed{\text{SIN}}(2x)$ in its fourth significant digit. The same identities are satisfied by $\boxed{\text{TRIG}}(x)$ values as by $\text{trig}(x)$ values even though $\boxed{\text{TRIG}}(x)$ and $\text{trig}(x)$ may disagree.

- Despite the two kinds of errors in $\boxed{\text{TRIG}}$, its computed values preserve familiar relationships wherever possible:

- Sign symmetry:
$$\begin{aligned} \boxed{\text{COS}}(-x) &= \boxed{\text{COS}}(x) \\ \boxed{\text{SIN}}(-x) &= -\boxed{\text{SIN}}(x) \end{aligned}$$

- **Monotonicity:** if $\text{trig}(x) \geq \text{trig}(y)$,
then $\boxed{\text{TRIG}}(x) \geq \boxed{\text{TRIG}}(y)$
(provided $|x - y| < 3$)
- **Limiting inequalities:** $\boxed{\text{SIN}}(x)/x \leq 1$ for all $x \neq 0$
 $\boxed{\text{TAN}}(x)/x \geq 1$ for $0 < |x| < \pi/2$
 $-1 \leq \boxed{\text{SIN}}(x)$ and $\boxed{\text{COS}}(x) \leq 1$
for all x

What do these properties imply for engineering calculations? *You don't have to remember them!*

In general, engineering calculations will not be affected by the difference between p and π , because the consequences of that difference in the formula defining $\boxed{\text{TRIG}}(x)$ above are swamped by the difference between $\boxed{\pi}$ and π and by ordinary unavoidable roundoff in x or in $\text{trig}(x)$. For engineering purposes, the ratio $\pi/p = 0.99999999999999342\dots$ could be replaced by 1 without visible effect upon the behavior of $\boxed{\text{TRIG}}$.

Example 5: Lunar Phases. If the distance between our Earth and its moon were known accurately, we could calculate the phase difference between radar signals transmitted to and reflected from the moon. In this calculation the phase shift introduced by $p \neq \pi$ has less effect than changing the distance between Earth and moon by as little as the thickness of this page. Moreover, the calculation of the strength, direction, and rate of change of radiated signals near the moon or reflected signals near the Earth, calculations that depend upon the trigonometric identities' continuing validity, are unaffected by the fact that $p \neq \pi$; they rely instead upon the fact that p is a constant (independent of x in the formula for $\boxed{\text{TRIG}}(x)$), and that constant is very near π .

The HP-15C's keyboard functions that involve p are the trigonometric functions $\boxed{\text{SIN}}$, $\boxed{\text{COS}}$, and $\boxed{\text{TAN}}$ for real and complex arguments; hyperbolic functions $\boxed{\text{SINH}}$, $\boxed{\text{COSH}}$, and $\boxed{\text{TANH}}$ for complex arguments; complex operations $\boxed{e^x}$, $\boxed{10^x}$, and $\boxed{y^x}$; and real and complex $\boxed{\rightarrow R}$.

It all seems like much ado about very little. After a blizzard of formulas and examples, we conclude that the error caused by $p \neq \pi$ is negligible for engineering purposes, so we need not have bothered to know about it. That is the burden that conscientious error analysts must bear; if they merely took for granted that small errors are negligible, they might be wrong.

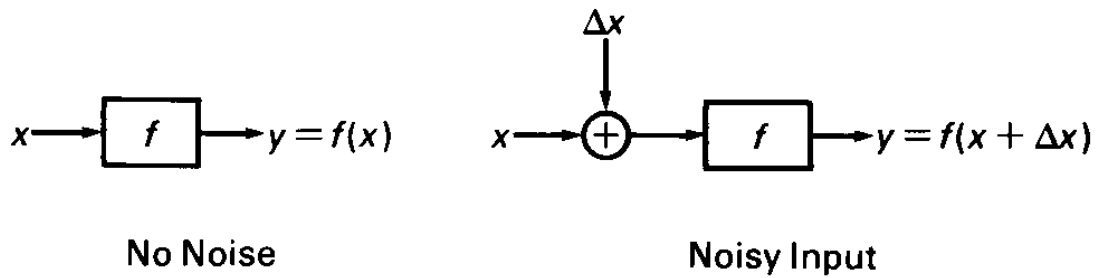
Backward Error Analysis

Until the late 1950's, most computer experts inclined to paranoia in their assessments of the damage done to numerical computations by rounding errors. To justify their paranoia, they could cite published error analyses like the one from which a famous scientist concluded that matrices as large as 40×40 were almost certainly impossible to invert numerically in the face of roundoff. However, by the mid 1960's matrices as large as 100×100 were being inverted routinely, and nowadays equations with hundreds of thousands of unknowns are being solved during geodetic calculations worldwide. How can we reconcile these accomplishments with the fact that that famous scientist's mathematical analysis was quite correct?

We understand better now than then why different formulas to calculate the same result might differ utterly in their degradation by rounding errors. For instance, we understand why the normal equations belonging to certain least-squares problems can be solved only in arithmetic carrying extravagantly high precision; this is what that famous scientist actually proved. We also know new procedures (one is presented on page 140) that can solve the same least-squares problems without carrying much more precision than suffices to represent the data. The new and better numerical procedures are not obvious, and might never have been found but for new and better techniques of error analysis by which we have learned to distinguish formulas that are hypersensitive to rounding errors from formulas that aren't. One of the new (in 1957) techniques is now called "backward error analysis," and you have already seen it in action twice: first, it explained why the procedure that calculates $\lambda(x)$ is accurate enough to dispel the inaccuracy in example 2; next, it explained why the calculator's `TRIG` functions very nearly satisfy the same identities as are satisfied by trig functions even for huge radian arguments x at which `TRIG`(x) and $\text{trig}(x)$ can be very different. The following paragraphs explain backward error analysis itself in general terms.

Consider some system F intended to transform an input x into an output $y = f(x)$. For instance, F could be a signal amplifier, a filter, a transducer, a control system, a refinery, a country's economy, a computer program, or a calculator. The input and output need not be numbers; they could be sets of numbers or matrices or anything else quantitative. Were the input x to be contaminated by noise Δx ,

then in consequence the output $y + \Delta y = f(x + \Delta x)$ would generally be contaminated by noise $\Delta y = f(x + \Delta x) - f(x)$.

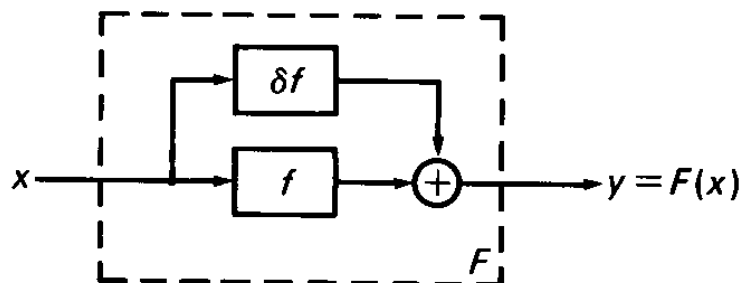


Some transformations f are stable in the presence of input noise; they keep Δy relatively small as long as Δx is relatively small. Other transformations f may be unstable in the presence of noise because certain relatively small input noises Δx cause relatively huge perturbations Δy in the output. In general, the input noise Δx will be colored in some way by the intended transformation f on the way from input to output noise Δy , and no diminution in Δy is possible without either diminishing Δx or changing f . Having accepted f as a specification for performance or as a goal for design, we must acquiesce to the way f colors noise at its input.

The real system F differs from the intended f because of noise or other discrepancies inside F . Before we can appraise the consequences of that internal noise we must find a way to represent it, a notation. The simplest way is to write

$$F(x) = (f + \delta f)(x)$$

where the perturbation δf represents the internal noise in F .



One Small Output Perturbation (Level 1)

We hope the noise term δf is negligible compared with f . When that hope is fulfilled, we classify F in Level 1 for the purposes of

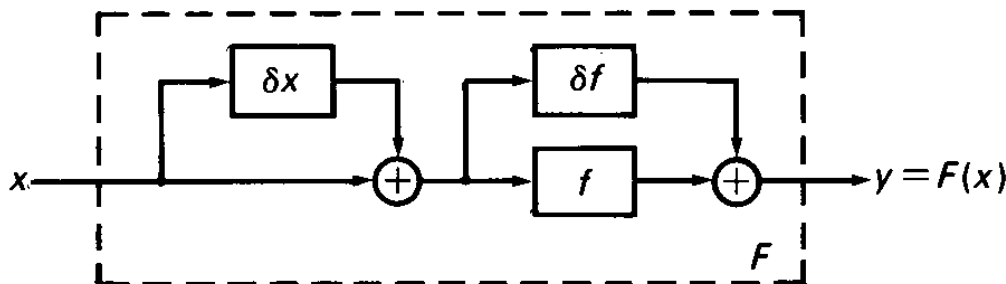
exposition; this means that the noise internal to F can be explained as *one* small addition δf to the intended output f .

For example, $F(x) = \boxed{\text{LN}}(x)$ is classified in Level 1 because the dozens of small errors committed by the HP-15C during its calculation of $F(x) = (f + \delta f)(x)$ amounts to a perturbation $\delta f(x)$ smaller than 0.6 in the last (10th) significant digit of the desired output $f(x) = \ln(x)$. But $F(x) = \boxed{\text{SIN}}(x)$ is not in Level 1 for radian x because $F(x)$ can differ too much from $f(x) = \sin(x)$; for instance $F(10^{14} \boxed{\pi}) = 0.799\dots$ is opposite in sign from $f(10^{14} \boxed{\pi}) = -0.784\dots$, so the equation $F(x) = (f + \delta f)(x)$ can be true only if δf is sometimes rather bigger than f , which looks bad.

Real systems more often resemble $\boxed{\text{SIN}}$ than $\boxed{\text{LN}}$. Noise in most real systems can accumulate occasionally to swamp the desired output, at least for some inputs, and yet such systems do not necessarily deserve condemnation. Many a real system F operates reliably because its internal noise, though sometimes large, never causes appreciably more harm than might be caused by some tolerably small perturbation δx to the input signal x . Such systems can be represented as

$$F(x) = (f + \delta f)(x + \delta x)$$

where δf is always small compared with f and δx is always smaller than or comparable with the noise Δx expected to contaminate x . The two noise terms δf and δx are hypothetical noises introduced to explain diverse noise sources actually distributed throughout F . Some of the noise appears as a tolerably small perturbation δx to the input—hence the term “backward error analysis.” Such a system F , whose noise can be accounted for by two tolerably small perturbations, is therefore classified into Level 2 for purposes of exposition.



Small Input and Output Perturbations (Level 2)

No difference will be perceived at first between Level 1 and Level 2 by readers accustomed to linear systems and small signals because such systems' errors can be referred indiscriminately to output or input. However, other more general systems that are digital or nonlinear do not admit arbitrary reattribution of output noise to input noise nor vice-versa.

For example, can all the error in $\boxed{\text{COS}}$ be attributed, merely by writing $\boxed{\text{COS}}(x) = \cos(x + \delta x)$, to an input perturbation δx small compared with the input x ? Not when x is very small. For instance, when x approaches 10^{-5} radians, then $\cos(x)$ falls very near 0.99999999995 and must then round to either $1 = \cos(0)$ or $0.9999999999 = \cos(1.414... \times 10^{-5})$. Therefore $\boxed{\text{COS}}(x) = \cos(x + \delta x)$ is true only if δx is allowed to be relatively large, nearly as large as x when x is very small. If we wish to explain the error in $\boxed{\text{COS}}$ by using only relatively small perturbations, we need at least two of them: one a perturbation $\delta x = (-6.58... \times 10^{-14})x$ smaller than roundoff in the input; and another in the output comparable with roundoff there, so that $\boxed{\text{COS}}(x) = (\cos + \delta \cos)(x + \delta x)$ for some unknown $|\delta \cos| \leq (6 \times 10^{-10})|\cos|$.

Like $\boxed{\text{COS}}$, every system F in Level 2 is characterized by just two small tolerances—call them ϵ and η —that sum up all you have to know about that system's internal noise. The tolerance ϵ constrains a hypothetical output noise, $|\delta f| \leq \epsilon|f|$, and η constrains a hypothetical input noise, $|\delta x| \leq \eta|x|$, that might appear in a simple formula like

$$F(x) = (f + \delta f)(x + \delta x) \quad \text{for } |\delta f| \leq \epsilon|f| \quad \text{and} \quad |\delta x| \leq \eta|x|.$$

The goal of backward error analysis is to ascertain that all the internal noise of F really can be encompassed by so simple a formula with satisfactorily small tolerances ϵ and η . At its best, backward error analysis confirms that *the realized value $F(x)$ scarcely differs from the ideal value $f(x + \delta x)$ that would have been produced by an input $x + \delta x$ scarcely different from the actual input x* , and gives the word "scarcely" a quantitative meaning (ϵ and η). But, backward error analysis succeeds only for systems F designed very carefully to ensure that every internal noise source is equivalent at worst to a tolerably small input or output perturbation. First attempts at system design, especially programs to perform numerical computations, often suffer from internal noise in a more complicated and disagreeable way illustrated by the following example.

Example 6: The Smaller Root of a Quadratic. The two roots x and y of the quadratic equation $c - 2bz + az^2 = 0$ are real whenever $d = b^2 - ac$ is nonnegative. Then the root y of smaller magnitude can be regarded as a function $y = f(a, b, c)$ of the quadratic's coefficients

$$f(a, b, c) = \begin{cases} (b - \sqrt{d} \operatorname{sgn}(b))/a & \text{if } a \neq 0 \\ (c/b)/2 & \text{otherwise.} \end{cases}$$

Were this formula translated directly in a program $F(a, b, c)$ intended to calculate $f(a, b, c)$, then whenever ac is so small compared with b^2 that the computed value of d rounds to b^2 , that program could deliver $F = 0$ even though $f \neq 0$. So drastic an error cannot be explained by backward error analysis because no relatively small perturbations to each coefficient a , b , and c could drive c to zero, as would be necessary to change the smaller root y into 0. On the other hand, the algebraically equivalent formula

$$f(a, b, c) = \begin{cases} c/(b + \sqrt{d} \operatorname{sgn}(b)) & \text{if divisor is nonzero} \\ 0 & \text{otherwise} \end{cases}$$

translates into a much more accurate program F whose errors do no more damage than would a perturbation in the last (10th) significant digit of c . Such a program will be listed later (page 205) and must be used in those instances, common in engineering, when the smaller root y is needed accurately despite the fact that the quadratic's other unwanted root is relatively large.

Almost all the functions built into the HP-15C have been designed so that backward error analysis will account for their errors satisfactorily. The exceptions are **SOLVE**, **f_i**, and the statistics keys **s**, **L.R.**, and **y_r** which can malfunction in certain pathological cases. Otherwise, every calculator function F intended to produce $f(x)$ produces instead a value $F(x)$ no farther from $f(x)$ than if first x had been perturbed to $x + \delta x$ with $|\delta x| \leq \eta|x|$, then $f(x + \delta x)$ were perturbed to $(f + \delta f)(x + \delta x)$ with $|\delta f| \leq \epsilon|f|$. The tolerances η and ϵ vary a little from function to function; roughly speaking,

$$\begin{aligned} \eta = 0 \text{ and } \epsilon < 10^{-9} & \quad \text{for all functions in Level 1,} \\ \eta < 10^{-12} \text{ and } \epsilon < 6 \times 10^{-10} & \quad \text{for other real and complex functions.} \end{aligned}$$

For matrix operations, the magnitudes $|\delta x|$, $|x|$, $|\delta f|$, and $|f|$ must be replaced by matrix norms $\|\delta \mathbf{x}\|$, $\|\mathbf{x}\|$, $\|\delta \mathbf{f}\|$, and $\|\mathbf{f}\|$ respectively, which are explained in section 4 and evaluated using [MATRIX] 7 or [MATRIX] 8. Then all matrix functions not in Level 1 fall into Level 2 with roughly

$$\begin{array}{ll} \eta \leq 10^{-12}n \text{ and } \epsilon < 10^{-9} & \text{for matrix operations (other than} \\ & \text{determinant [MATRIX] 9, } \boxed{\div}, \text{ and } \boxed{1/x}) \\ \eta \leq 10^{-9}n \text{ and } \epsilon < 10^{-9} & \text{for determinant [MATRIX] 9, } \boxed{1/x}, \\ & \text{and } \boxed{\div} \text{ with a matrix divisor} \end{array}$$

where n is the largest dimension of any matrix involved in the operation.

The implications of successful backward error analysis look simple only when the input data x comes contaminated by unavoidable and uncorrelated noise Δx , as is often the case. Then when we wish to calculate $f(x)$, the best we could hope to get is $f(x + \Delta x)$, but we actually get $F(x + \Delta x) = (f + \delta f)(x + \Delta x + \delta x)$, where $|\delta f| \leq \epsilon|f|$ and $|\delta x| \leq \eta|x|$.

What we get is scarcely worse than the best we could hope for provided the tolerances ϵ and η are small enough, particularly if $|\Delta x|$ is likely to be at least roughly as big as $\eta|x|$. Of course, the best we could hope for may be very bad, especially if f possesses a singularity closer to x than the tolerances upon x 's perturbations Δx and δx .

Backward Error Analysis Versus Singularities

The word "singularity" refers to both a special value of the argument x and to the way $f(x)$ misbehaves as x approaches that special value. Most commonly, $f(x)$ or its first derivative $f'(x)$ may become infinite or violently oscillatory as x approaches the singularity. Sometimes the singularities of $\ln|f|$ are called singularities of f , thereby including the zeros of f among its singularities; this makes sense when the relative accuracy of a computation of f is at issue, as we shall see. For our purposes the meaning of "singularity" can be left a little vague.

What we usually want to do with singularities is avoid or neutralize them. For instance, the function

$$c(x) = \begin{cases} (1 - \cos x)/x^2 & \text{if } x \neq 0 \\ 1/2 & \text{otherwise} \end{cases}$$

has no singularity at $x = 0$ even though its constituents $1 - \cos x$ and x^2 (actually, their logarithms) do behave singularly as x approaches 0. The constituent singularities cause trouble for the program that calculates $c(x)$. Most of the trouble is neutralized by the choice of a better formula

$$c(x) = \begin{cases} \frac{1}{2} \left(\frac{\sin(x/2)}{x/2} \right)^2 & \text{if } x/2 \neq 0 \\ 1/2 & \text{otherwise.} \end{cases}$$

Now the singularity can be avoided entirely by testing whether $x/2 = 0$ in the program that calculates $c(x)$.

Backward error analysis complicates singularities in a way that is easiest to illustrate with the function $\lambda(x) = \ln(1 + x)$ that solved the savings problem in example 2. The procedure used there calculated $u = 1 + x$ (rounded) $= 1 + x + \Delta x$. Then

$$\lambda(x) = \begin{cases} x & \text{if } u = 1 \\ \ln(u) x / (u - 1) & \text{otherwise.} \end{cases}$$

This procedure exploits the fact that $\lambda(x)/x$ has a removable singularity at $x = 0$, which means that $\lambda(x)/x$ varies continuously and approaches 1 as x approaches 0. Therefore, $\lambda(x)/x$ is relatively closely approximated by $\lambda(x + \Delta x)/(x + \Delta x)$ when $|\Delta x| < 10^{-9}$, and hence

$$\lambda(x) = x(\lambda(x)/x) \approx x(\lambda(x + \Delta x)/(x + \Delta x)) = x(\ln(u)/(u - 1)),$$

all calculated accurately because $\boxed{\text{LN}}$ is in Level 1. What might happen if $\boxed{\text{LN}}$ were in Level 2 instead?

If $\boxed{\text{LN}}$ were in Level 2, then “successful” backward error analysis would show that, for arguments u near 1, $\boxed{\text{LN}}(u) = \ln(u + \delta u)$ with $|\delta u| < 10^{-9}$. Then the procedure above would produce not $x(\ln(u)/(u - 1))$, but

$$\begin{aligned} x(\ln(u + \delta u)/(u - 1)) &= x\lambda(x + \Delta x + \delta u)/(x + \Delta x) \\ &= x(\lambda(x + \Delta x + \delta u)/(x + \Delta x + \delta u)) \frac{x + \Delta x + \delta u}{x + \Delta x} \\ &\approx x(\lambda(x)/x)(1 + \delta u/(x + \Delta x)) \\ &= \lambda(x)(1 + \delta u/(x + \Delta x)). \end{aligned}$$

When $|x + \Delta x|$ is not much bigger than 10^{-9} , the last expression can be utterly different from $\lambda(x)$. Therefore, the procedure that solved example 2 would fail on machines whose $\boxed{\text{LN}}$ is not in Level 1. There *are* such machines, and on them the procedure does collapse for certain otherwise innocuous inputs. Similar failures also occur on machines that produce $(u + \delta' u) - 1$ instead of $u - 1$ because their $\boxed{-}$ function lies in Level 2 instead of Level 1. And those machines that produce $\ln(u + \delta u)/(u + \delta' u - 1)$ instead of $\ln(u)/(u - 1)$, because both $\boxed{\text{LN}}$ and $\boxed{-}$ lie in Level 2, would be doubly vulnerable but for an ill-understood accident that usually correlates the two backward errors δu and $\delta' u$ in such a way as causes only half the significant digits of the computed λ , instead of all of them, to be wrong.

Summary to Here

Now that the complexity injected by backward error analysis into singularities has been exposed, the time has come to summarize, to simplify, and to consolidate what has been discussed so far.

- Many numerical procedures produce results too wrong to be justified by any satisfactory error analysis, backward or not.
- Some numerical procedures produce results only slightly worse than would have been obtained by exactly solving a problem differing only slightly from the given problem. Such procedures, classified in Level 2 for our purposes, are widely accepted as satisfactory from the point of view of backward error analysis.
- Procedures in Level 2 can produce results relatively far from what would have been obtained had no errors at all been committed, but large errors can result only for data relatively near a singularity of the function being computed.
- Procedures in Level 1 produce relatively accurate results regardless of near approach to a singularity. Such procedures are rare, but preferable if only because their results are easier to interpret, especially when several variables are involved.

A simple example illustrates all four points.

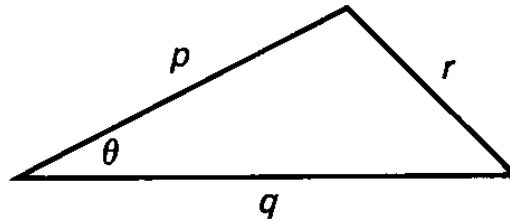
Example 7: The Angle in a Triangle. The cosine law for triangles says

$$r^2 = p^2 + q^2 - 2pq \cos \theta$$

for the figure shown below. Engineering and scientific calculations often require that the angle θ be calculated from given values p , q , and r for the length of the triangle's sides. This calculation is feasible provided $0 < p \leq q + r$, $0 < q \leq p + r$, and $0 \leq r \leq p + q$, and then

$$0 \leq \theta = \cos^{-1}(((p^2 + q^2) - r^2)/(2pq)) \leq 180^\circ;$$

otherwise, no triangle exists with those side lengths, or else $\theta = 0/0$ is indeterminate.



The foregoing formula for θ defines a function $\theta = f(p, q, r)$ and also in a natural way, a program $F(p, q, r)$ intended to calculate the function. That program is labeled "A" below, with results $F_A(p, q, r)$ tabulated for certain inputs p , q , and r corresponding to sliver-shaped triangles for which the formula suffers badly from roundoff. The numerical unreliability of this formula is well known as is that of the algebraically equivalent but more reliable formula $\theta = f(p, q, r) = 2 \tan^{-1} \sqrt{ab/(cs)}$, where $s = (p + q + r)/2$, $a = s - p$, $b = s - q$, and $c = s - r$. Another program $F(p, q, r)$ based upon this better formula is labeled "B" below, with results $F_B(p, q, r)$ for selected inputs. Apparently F_B is not much more reliable than F_A . Most of the poor results could be explained by backward error analysis if we assume that the calculations yield $F(p, q, r) = f(p + \delta p, q + \delta q, r + \delta r)$ for unknown but small perturbations satisfying $|\delta p| < 10^{-9}|p|$, etc. Even if this explanation were true, it would have perplexing and disagreeable consequences, because the angles in sliver-shaped triangles can change relatively drastically when the sides are perturbed relatively slightly; $f(p, q, r)$ is relatively unstable for marginal inputs.

Actually the preceding explanation is false. No backward error analysis could account for the results tabulated for F_A and F_B under case 1 below unless perturbations δp , δq , and δr were allowed to corrupt the fifth significant digit of the input, changing 1 to 1.0001 or 0.9999. That much is too much noise to tolerate in a 10-digit calculation. A better program by far is F_C , labeled "C" and explained shortly afterwards.

The three bottom lines in the table below show results for three programs "A", "B", and "C" based upon three different formulas $F(p, q, r)$ all algebraically equivalent to

$$\theta = f(p, q, r) = \cos^{-1}((p^2 + q^2 - r^2)/(2pq)).$$

Disparate Results from Three Programs F_A, F_B, F_C

	Case 1	Case 2	Case 3
p	1.	9.999999996	10.
q	1.	9.999999994	5.000000001
r	1.00005×10^{-5}	3×10^{-9}	15.
F_A	0.	0.	180.
F_B	5.73072×10^{-4}	Error 0	180.
F_C	5.72986×10^{-4}	1.28117×10^{-8}	179.9985965
	Case 4	Case 5	Case 6
p	0.527864055	9.999999996	9.999999999
q	9.472135941	3×10^{-9}	9.999999999
r	9.999999996	9.999999994	20.
F_A	Error 0	48.18968509	180.
F_B	Error 0	Error 0	180.
F_C	180.	48.18968510	Error 0
	Case 7	Case 8	Case 9
p	1.00002	3.162277662	3.162277662
q	1.00002	2.3×10^{-9}	1.5555×10^{-6}
r	2.00004	3.162277661	3.162277661
F_A	Error 0	90.	90.
F_B	180.	70.52877936	89.96318706
F_C	180.	64.22853822	89.96315156

To use a program, key in p q r , run program "A", "B", or "C", and wait to see the program's approximation F to $\theta = f$. Only program "C" is reliable.

Keystrokes	Display
\boxed{g} \boxed{DEG}	
\boxed{g} $\boxed{P/R}$	
\boxed{f} \boxed{CLEAR} \boxed{PRGM}	000-
\boxed{f} \boxed{LBL} \boxed{A}	001-42,21,11
\boxed{g} $\boxed{x^2}$	002- 43 11
$\boxed{x} \boxed{\geq} \boxed{y}$	003- 34
\boxed{g} $\boxed{x^2}$	004- 43 11
\boxed{g} \boxed{LSTx}	005- 43 36
\boxed{g} $\boxed{R\uparrow}$	006- 43 33
\boxed{x}	007- 20
$\boxed{x} \boxed{\geq} \boxed{y}$	008- 34
\boxed{g} \boxed{LSTx}	009- 43 36
\boxed{g} $\boxed{x^2}$	010- 43 11
$\boxed{+}$	011- 40
\boxed{g} $\boxed{R\uparrow}$	012- 43 33
$\boxed{-}$	013- 30
$\boxed{x} \boxed{\geq} \boxed{y}$	014- 34
\boxed{ENTER}	015- 36
$\boxed{+}$	016- 40
$\boxed{\div}$	017- 10
\boxed{g} $\boxed{COS^{-1}}$	018- 43 24
\boxed{g} \boxed{RTN}	019- 43 32
\boxed{f} \boxed{LBL} \boxed{B}	020-42,21,12
\boxed{STO} $\boxed{1}$	021- 44 1
\boxed{ENTER}	022- 36
\boxed{g} $\boxed{R\uparrow}$	023- 43 33
\boxed{STO} $\boxed{+}$ $\boxed{1}$	024-44,40, 1
\boxed{g} $\boxed{R\uparrow}$	025- 43 33
\boxed{STO} $\boxed{+}$ $\boxed{1}$	026-44,40, 1
$\boxed{2}$	027- 2
\boxed{STO} $\boxed{\div}$ $\boxed{1}$	028-44,10, 1
$\boxed{R\downarrow}$	029- 33
\boxed{RCL} $\boxed{-}$ $\boxed{1}$	030-45,30, 1
$\boxed{x} \boxed{\geq} \boxed{y}$	031- 34
\boxed{RCL} $\boxed{-}$ $\boxed{1}$	032-45,30, 1
\boxed{x}	033- 20
$\boxed{\sqrt{x}}$	034- 11
$\boxed{x} \boxed{\geq} \boxed{y}$	035- 34
\boxed{RCL} $\boxed{-}$ $\boxed{1}$	036-45,30, 1
\boxed{RCL} $\boxed{\times}$ $\boxed{1}$	037-45,20, 1

Keystrokes	Display
CHS	038- 16
\sqrt{x}	039- 11
g \rightarrow P	040- 43 1
R \downarrow	041- 33
x	042- 20
g RTN	043- 43 32
f LBL C	044-42,21,13
STO 0	045- 44 0
R \downarrow	046- 33
g $x \leq y$	047- 43 10
$x \geq y$	048- 34
STO 1	049- 44 1
STO + 0	050-44,40, 0
$x \geq y$	051- 34
STO + 0	052-44,40, 0
-	053- 30
g R \uparrow	054- 43 33
STO - 1	055-44,30, 1
g LSTx	056- 43 36
ENTER	057- 36
RCL + 1	058-45,40, 1
\sqrt{x}	059- 11
f $x \geq 0$	060-42, 4, 0
\sqrt{x}	061- 11
STO x 0	062-44,20, 0
g CLx	063- 43 35
+	064- 40
R \downarrow	065- 33
+	066- 40
f $x \geq 1$	067-42, 4, 1
g R \uparrow	068- 43 33
g LSTx	069- 43 36
g $x \leq y$	070- 43 10
GTO .9	071- 22 .9
R \downarrow	072- 33
g TEST 2	073-43,30, 2
\sqrt{x}	074- 11
$x \geq y$	075- 34
GTO .8	076- 22 .8
f LBL .9	077-42,21, .9

Keystrokes	Display
\boxed{g} $\boxed{\text{TEST}}$ $\boxed{2}$	078-43,30, 2
$\boxed{\sqrt{x}}$	079- 11
\boxed{g} $\boxed{R\uparrow}$	080- 43 33
\boxed{f} $\boxed{\text{LBL}}$ $\boxed{.8}$	081-42,21, .8
$\boxed{-}$	082- 30
$\boxed{\sqrt{x}}$	083- 11
$\boxed{\text{RCL}}$ $\boxed{1}$	084- 45 1
$\boxed{\sqrt{x}}$	085- 11
$\boxed{\times}$	086- 20
$\boxed{\text{RCL}}$ $\boxed{0}$	087- 45 0
\boxed{g} $\boxed{\rightarrow P}$	088- 43 1
\boxed{g} $\boxed{x=0}$	089- 43 20
$\boxed{\div}$	090- 10
$\boxed{x \geq y}$	091- 34
$\boxed{\text{ENTER}}$	092- 36
$\boxed{+}$	093- 40
\boxed{g} $\boxed{\text{RTN}}$	094- 43 32
\boxed{g} $\boxed{\text{P/R}}$	

The results $F_C(p, q, r)$ are correct to at least nine significant digits. They are obtained from a program "C" that is utterly reliable though rather longer than the unreliable programs "A" and "B". The method underlying program "C" is:

1. If $p < q$, then swap them to ensure $p \geq q$.
2. Calculate $b = (p - q) + r$, $c = (p - r) + q$, and $s = (p + r) + q$.
3. Calculate

$$a = \begin{cases} r - (p - q) & \text{if } q \geq r \geq 0 \\ q - (p - r) & \text{if } r > q \geq 0 \\ \text{Error 0} & \text{otherwise (no triangle exists).} \end{cases}$$

4. Calculate $F_C(p, q, r) = 2 \tan^{-1}(\sqrt{ab}/\sqrt{cs})$.

This procedure delivers $F_C(p, q, r) = \theta$ correct to almost nine significant digits, a result surely easier to use and interpret than the results given by the other better-known formulas. But this procedure's internal workings are hard to explain; indeed, the procedure may malfunction on some calculators and computers.

The procedure works impeccably on only certain machines like the HP-15C, whose subtraction operation is free from avoidable error and therefore enjoys the following property: Whenever y lies between $x/2$ and $2x$, the subtraction operation introduces no roundoff error into the calculated value of $x - y$. Consequently, whenever cancellation might leave relatively large errors contaminating a , b , or c , the pertinent difference $(p - q)$ or $(p - r)$ turns out to be free from error, and then cancellation turns out to be advantageous!

Cancellation remains troublesome on those other machines that calculate $(x + \delta x) - (y + \delta y)$ instead of $x - y$ even though neither δx nor δy amounts to as much as one unit in the last significant digit carried in x or y respectively. Those machines deliver $F_C(p, q, r) = f(p + \delta p, q + \delta q, r + \delta r)$ with end-figure perturbations δp , δq , and δr that always seem negligible from the viewpoint of backward error analysis, but which can have disconcerting consequences. For instance, only one of the triples (p, q, r) or $(p + \delta p, q + \delta q, r + \delta r)$, not both, might constitute the edge lengths of a feasible triangle, so F_C might produce an error message when it shouldn't, or vice-versa, on those machines.

Backward Error Analysis of Matrix Inversion

The usual measure of the magnitude of a matrix \mathbf{X} is a norm $\|\mathbf{X}\|$ such as is calculated by either **MATRIX** 7 or **MATRIX** 8; we shall use the former norm, the row norm

$$\|\mathbf{X}\| = \max_i \sum_j |x_{ij}|$$

in what follows. This norm has properties similar to those of the length of a vector and also the multiplicative property

$$\|\mathbf{XY}\| \leq \|\mathbf{X}\| \|\mathbf{Y}\|.$$

When the equation $\mathbf{Ax} = \mathbf{b}$ is solved numerically with a given $n \times n$ matrix \mathbf{A} and column vector \mathbf{b} , the calculated solution is a column vector \mathbf{c} which satisfies nearly the same equation as does \mathbf{x} , namely

$$(\mathbf{A} + \delta\mathbf{A})\mathbf{c} = \mathbf{b}$$

with $\|\delta\mathbf{A}\| < 10^{-9} n \|\mathbf{A}\|$.

Consequently the residual $\mathbf{b} - \mathbf{A}\mathbf{c} = (\delta\mathbf{A})\mathbf{c}$ is always relatively small; quite often the residual norm $\|\mathbf{b} - \mathbf{A}\mathbf{c}\|$ is smaller than $\|\mathbf{b} - \mathbf{A}\bar{\mathbf{x}}\|$ where $\bar{\mathbf{x}}$ is obtained from the true solution \mathbf{x} by rounding each of its elements to 10 significant digits. Consequently, \mathbf{c} can differ significantly from \mathbf{x} only if \mathbf{A} is nearly singular, or equivalently only if $\|\mathbf{A}^{-1}\|$ is relatively large compared with $1/\|\mathbf{A}\|$;

$$\begin{aligned} \|\mathbf{x} - \mathbf{c}\| &= \|\mathbf{A}^{-1}(\mathbf{b} - \mathbf{A}\mathbf{c})\| \\ &\leq \|\mathbf{A}^{-1}\| \|\delta\mathbf{A}\| \|\mathbf{c}\| \\ &\leq 10^{-9} n \|\mathbf{c}\| / \sigma(\mathbf{A}) \end{aligned}$$

where $\sigma(\mathbf{A}) = 1/(\|\mathbf{A}\| \|\mathbf{A}^{-1}\|)$ is the reciprocal of the condition number and measures how relatively near to \mathbf{A} is the nearest singular matrix \mathbf{S} , since

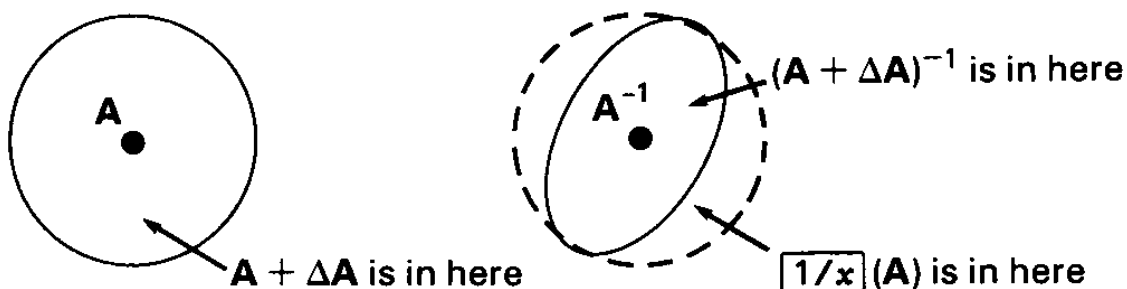
$$\min_{\det(\mathbf{S})=0} \|\mathbf{A} - \mathbf{S}\| = \sigma(\mathbf{A}) \|\mathbf{A}\|.$$

These relations and some of their consequences are discussed extensively in section 4.

The calculation of \mathbf{A}^{-1} is more complicated. Each column of the calculated inverse $\boxed{1/\mathbf{x}}(\mathbf{A})$ is the corresponding column of some $(\mathbf{A} + \delta\mathbf{A})^{-1}$, but each column has its own small $\delta\mathbf{A}$. Consequently, no single small $\delta\mathbf{A}$, with $\|\delta\mathbf{A}\| \leq 10^{-9} n \|\mathbf{A}\|$, need exist satisfying

$$\|(\mathbf{A} + \delta\mathbf{A})^{-1} - \boxed{1/\mathbf{x}}(\mathbf{A})\| \leq 10^{-9} \|\boxed{1/\mathbf{x}}(\mathbf{A})\|$$

roughly. Usually such a $\delta\mathbf{A}$ exists, but not always. This does not violate the prior assertion that the matrix operations $\boxed{1/\mathbf{x}}$ and $\boxed{\div}$ lie in Level 2; they are covered by the second assertion of the summary on page 194. The accuracy of $\boxed{1/\mathbf{x}}(\mathbf{A})$ can be described in terms of the inverses of all matrices $\mathbf{A} + \Delta\mathbf{A}$ so near \mathbf{A} that $\|\Delta\mathbf{A}\| \leq 10^{-9} n \|\mathbf{A}\|$; the worst among those $(\mathbf{A} + \Delta\mathbf{A})^{-1}$ is at least about as far from \mathbf{A}^{-1} in norm as the calculated $\boxed{1/\mathbf{x}}(\mathbf{A})$. The figure below illustrates the situation.



As $\mathbf{A} + \Delta\mathbf{A}$ runs through matrices with $\|\Delta\mathbf{A}\|$ at least about as large as roundoff in $\|\mathbf{A}\|$, its inverse $(\mathbf{A} + \Delta\mathbf{A})^{-1}$ must roam at least about as far from \mathbf{A}^{-1} as the distance from \mathbf{A}^{-1} to the computed $\boxed{1/x}(\mathbf{A})$. All these excursions are very small unless \mathbf{A} is too near a singular matrix, in which case the matrix should be preconditioned away from near singularity. (Refer to section 4.)

If among those neighboring matrices $\mathbf{A} + \Delta\mathbf{A}$ lurk some that are singular, then many $(\mathbf{A} + \Delta\mathbf{A})^{-1}$ and $\boxed{1/x}(\mathbf{A})$ may differ utterly from \mathbf{A}^{-1} . However, the residual norm will always be relatively small:

$$\frac{\|\mathbf{A}(\mathbf{A} + \Delta\mathbf{A})^{-1} - \mathbf{I}\|}{\|\mathbf{A}\| \|(\mathbf{A} + \Delta\mathbf{A})^{-1}\|} \leq \frac{\|\Delta\mathbf{A}\|}{\|\mathbf{A}\|} \leq 10^{-9}n.$$

This last inequality remains true when $\boxed{1/x}(\mathbf{A})$ replaces $(\mathbf{A} + \Delta\mathbf{A})^{-1}$.

If \mathbf{A} is far enough from singularity that all

$$1/\|(\mathbf{A} + \Delta\mathbf{A})^{-1}\| > 10^{-9}n\|\mathbf{A}\| \geq \|\Delta\mathbf{A}\|,$$

then also

$$\begin{aligned} \frac{\|\mathbf{A}^{-1} - (\mathbf{A} + \Delta\mathbf{A})^{-1}\|}{\|(\mathbf{A} + \Delta\mathbf{A})^{-1}\|} &\leq \frac{\|\Delta\mathbf{A}\| \|(\mathbf{A} + \Delta\mathbf{A})^{-1}\|}{1 - \|\Delta\mathbf{A}\| \|(\mathbf{A} + \Delta\mathbf{A})^{-1}\|} \\ &\leq \frac{10^{-9}n\|\mathbf{A}\| \|(\mathbf{A} + \Delta\mathbf{A})^{-1}\|}{1 - 10^{-9}n\|\mathbf{A}\| \|(\mathbf{A} + \Delta\mathbf{A})^{-1}\|}. \end{aligned}$$

This inequality also remains true when $\boxed{1/x}(\mathbf{A})$ replaces $(\mathbf{A} + \Delta\mathbf{A})^{-1}$, and then everything on the right-hand side can be calculated, so the error in $\boxed{1/x}(\mathbf{A})$ cannot exceed a knowable amount. In other words, the radius of the dashed ball in the figure above can be calculated.

The estimates above tend to be pessimistic. However, to show why nothing much better is true in general, consider the matrix

$$\mathbf{X} = \begin{bmatrix} 0.00002 & -50,000 & 50,000.03 & -45 \\ 0 & 50,000 & -50,000.03 & 45 \\ 0 & 0 & 0.00002 & -50,000.03 \\ 0 & 0 & 0 & 52,000 \end{bmatrix}$$

and

$$\mathbf{X}^{-1} = \begin{bmatrix} 50,000 & 50,000 & & p & q \\ 0 & 0.00002 & 50,000.03 & 48,076.98077... \\ 0 & 0 & 50,000 & 48,076.95192... \\ 0 & 0 & 0 & 0 & 0.00001923076923... \end{bmatrix}.$$

Ideally, $p = q = 0$, but the HP-15C's approximation to \mathbf{X}^{-1} , namely $\boxed{1/x}(\mathbf{X})$, has $q = 9,643.269231$ instead, a relative error

$$\frac{\|\mathbf{X}^{-1} - \boxed{1/x}(\mathbf{X})\|}{\|\mathbf{X}^{-1}\|} = 0.0964\dots,$$

nearly 10 percent. On the other hand, if $\mathbf{X} + \Delta\mathbf{X}$ differs from \mathbf{X} only in its second column where $-50,000$ and $50,000$ are replaced respectively by $-50,000.000002$ and $49,999.999998$ (altered in the 11th significant digit), then $(\mathbf{X} + \Delta\mathbf{X})^{-1}$ differs significantly from \mathbf{X}^{-1} only insofar as $p = 0$ and $q = 0$ must be replaced by $p = 10,000.00600\dots$ and $q = 9,615.396154\dots$. Hence,

$$\frac{\|\mathbf{X}^{-1} - (\mathbf{X} + \Delta\mathbf{X})^{-1}\|}{\|\mathbf{X}^{-1}\|} = 0.196\dots;$$

the relative error in $(\mathbf{X} + \Delta\mathbf{X})^{-1}$ is nearly twice that in $\boxed{1/x}(\mathbf{X})$. Do not try to calculate $(\mathbf{X} + \Delta\mathbf{X})^{-1}$ directly, but use instead the formula

$$(\mathbf{X} - \mathbf{c}\mathbf{b}^T)^{-1} = \mathbf{X}^{-1} + \mathbf{X}^{-1}\mathbf{c}\mathbf{b}^T\mathbf{X}^{-1} / (1 - \mathbf{b}^T\mathbf{X}^{-1}\mathbf{c}),$$

which is valid for any column vector \mathbf{c} and row vector \mathbf{b}^T , and specifically for

$$\mathbf{c} = \begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \end{bmatrix} \text{ and } \mathbf{b}^T = [0 \quad 0.000002 \quad 0 \quad 0].$$

Despite that

$$\|\mathbf{X}^{-1} - \boxed{1/x}(\mathbf{X})\| < \|\mathbf{X}^{-1} - (\mathbf{X} + \Delta\mathbf{X})^{-1}\|,$$

it can be shown that no very small end-figure perturbation $\delta\mathbf{X}$ exists for which $(\mathbf{X} + \delta\mathbf{X})^{-1}$ matches $\boxed{1/x}(\mathbf{X})$ to more than five significant digits in norm.

Of course, none of these horrible things could happen if \mathbf{X} were not so nearly singular. Because $\|\mathbf{X}\| \|\mathbf{X}^{-1}\| > 10^{10}$, a change in \mathbf{X} amounting to less than one unit in the 10th significant digit of $\|\mathbf{X}\|$ could make \mathbf{X} singular; such a change might replace one of the diagonal elements 0.00002 of \mathbf{X} by zero. Since \mathbf{X} is so nearly singular, the accuracy of $\boxed{1/\mathbf{x}}(\mathbf{X})$ in this case rather exceeds what might be expected in general. What makes this example special is bad scaling; \mathbf{X} was obtained from an unexceptional matrix

$$\tilde{\mathbf{X}} = \begin{bmatrix} 2. & -5. & 5.000003 & -4.5 \times 10^{-12} \\ 0 & 5. & -5.000003 & 4.5 \times 10^{-12} \\ 0 & 0 & 2. & -5.000003 \\ 0 & 0 & 0 & 5.2 \end{bmatrix}$$

by multiplying each row and each column by a carefully chosen power of 10. Compensatory division of the columns and rows of the equally unexceptional matrix

$$\tilde{\mathbf{X}}^{-1} = \begin{bmatrix} 0.5 & 0.5 & p & q \\ 0 & 0.2 & 0.5000003 & 0.4807698077\dots \\ 0 & 0 & 0.5 & 0.4807695192\dots \\ 0 & 0 & 0 & 0.1923076923\dots \end{bmatrix}$$

yielded \mathbf{X}^{-1} , with $p = q = 0$. The HP-15C calculates $\boxed{1/\mathbf{x}}(\tilde{\mathbf{X}}) = \tilde{\mathbf{X}}^{-1}$ except that $q = 0$ is replaced by $q = 9.6 \times 10^{-11}$, a negligible change. This illustrates how drastically the perceived quality of computed results can be altered by scaling. (Refer to section 4 for more information about scaling.)

Is Backward Error Analysis a Good Idea?

The only good thing to be said for backward error analysis is that it explains internal errors in a way that liberates a system's user from having to know about internal details of the system. Given two tolerances, one upon the input noise δx and one upon the output noise δf , the user can analyze the consequences of internal noise in

$$F(x) = (f + \delta f)(x + \delta x)$$

by studying the noise propagation properties of the ideal system f without further reference to the possibly complex internal structure of F .

But backward error analysis is no panacea; it may explain errors but not excuse them. Because it complicates computations involving singularities, we have tried to eliminate the need for it wherever we could. If we knew how to eliminate the need for backward error analysis from every function built into the calculator, and to do so at tolerable cost, we would do that and simplify life for everyone. That simplicity would cost too much speed and memory for today's technology. The next example will illustrate the trade-offs involved.

Example 6 Continued. The program listed below solves the real quadratic equation $c - 2bz + az^2 = 0$ for real or complex roots.

To use the program, key the real constants into the stack (c **ENTER** b **ENTER** a) and run program "A".

The roots x and y will appear in the X- and Y-registers. If the roots are complex, the **C** annunciator turns on, indicating that Complex mode has been activated. The program uses labels "A" and ".9" and the Index register (but none of the other registers 0 to .9); therefore, the program may readily be called as a subroutine by other programs. The calling programs (after clearing flag 8 if necessary) can discover whether roots are real or complex by testing flag 8, which gets set only if roots are complex.

The roots x and y are so ordered that $|x| \geq |y|$ except possibly when $|x|$ and $|y|$ agree to more than nine significant digits. The roots are as accurate as if the coefficient c had first been perturbed in its 10th significant digit, the perturbed equation had been solved exactly, and its roots rounded to 10 significant digits. Consequently, the computed roots match the given quadratic's roots to at least five significant digits. More generally, if the roots x and y agree to n significant digits for some positive $n \leq 5$, then they are correct to at least $10 - n$ significant digits unless overflow or underflow occurs.

Keystrokes	Display
g P/R	
f CLEAR PRGM	000-
f LBL A	001-42,21,11
ENTER	002- 36
g R↑	003- 43 33
x	004- 20
g LSTx	005- 43 36

Keystrokes	Display
$x \geq y$	006- 34
g $R \uparrow$	007- 43 33
STO I	008- 44 25
g x^2	009- 43 11
$-$	010- 30
g $TEST$ 1	011-43,30, 1
GTO .9	012- 22 .9
CHS	013- 16
\sqrt{x}	014- 11
f $x \geq I$	015-42, 4,25
g $TEST$ 2	016-43,30, 2
RCL $-$ I	017-45,30,25
g $TEST$ 3	018-43,30, 3
RCL $+$ I	019-45,40,25
g $TEST$ 0	020-43,30, 0
\div	021- 10
g $LSTx$	022- 43 36
g $R \uparrow$	023- 43 33
\div	024- 10
g RTN	025- 43 32
f LBL .9	026-42,21, .9
\sqrt{x}	027- 11
RCL I	028- 45 25
g $R \uparrow$	029- 43 33
\div	030- 10
$x \geq y$	031- 34
g $LSTx$	032- 43 36
\div	033- 10
f I	034- 42 25
$ENTER$	035- 36
f $Re \geq Im$	036- 42 30
CHS	037- 16
f $Re \geq Im$	038- 42 30
g RTN	039- 43 32
g P/R	

The method uses $d = b^2 - ac$.

If $d < 0$, then the roots are a complex conjugate pair

$$(b/a) \pm i\sqrt{-d}/a.$$

If $d \geq 0$, then the roots are real numbers x and y calculated by

$$s = b + \sqrt{d} \operatorname{sgn}(b)$$

$$x = s/a$$

$$y = \begin{cases} c/s & \text{if } s \neq 0 \\ 0 & \text{if } s = 0. \end{cases}$$

The s calculation avoids destructive cancellation.

When $a = 0 \neq b$, the larger root x , which should be ∞ , encounters division by zero (**Error 0**) that can be cleared by pressing **R↓** three times to exhibit the smaller root y correctly calculated. But when all three coefficients vanish, the **Error 0** message signals that both roots are arbitrary.

The results of several cases are summarized below.

	Case 1	Case 2	Case 3	Case 4
c	3	4	1	654,321
b	2	0	1	654,322
a	1	1	10^{-13}	654,323
Roots	Real	Complex	Real	Real
	3	$0 \pm 2i$	2×10^{13}	0.9999984717
	1		0.5	0.9999984717
	Case 5	Case 6		
c	46,152,709	12,066,163		
b	735,246	987,644		
a	11,713	80,841		
Roots	Real	Complex		
	62.77179203	$12.21711755 \pm i0.001377461$		
	62.77179203			

The last three cases show how severe are the results of perturbing the 10th significant digit of any coefficient of any quadratic whose roots are nearly coincident. The correct roots for these cases are

Case 4: 1 and 0.9999969434

Case 5: $62.77179203 \pm i8.5375 \times 10^{-5}$

Case 6: $12.21711755 \pm i0.001374514$

Despite errors in the fifth significant digit of the results, subroutine "A" suffices for almost all engineering and scientific applications of quadratic equations. Its results are correct to nine significant digits for most data, including c , b , and a representable exactly using only five significant digits; and the computed roots are correct to at least five significant digits in any case because they cannot be appreciably worse than if the data had been entered with errors in the 10th significant digit. Nonetheless, some readers will feel uneasy about results calculated to 10 significant digits but correct to only 5. If only to simplify their understanding of the relationship between input data and output results, they might still prefer roots correct to nine significant digits in all cases.

Programs do exist which, while carrying only 10 significant digits during arithmetic, will calculate the roots of any quadratic correctly to at least nine significant digits regardless of how nearly coincident those roots may be. All such programs calculate $d = b^2 - ac$ by some trick tantamount to carrying 20 significant digits whenever b^2 and ac nearly cancel, so those programs are a lot longer and slower than the simple subroutine "A" provided above. Subroutine "B" below, which uses such a trick,* is a very short program that guarantees nine correct significant digits on a 10-digit calculator. It uses labels "B", ".7", and ".8" and registers R_0 through R_9 and the Index register. To use it, key in c **ENTER** b **ENTER** a , run subroutine "B", and wait for results as before.

Keystrokes	Display
g P/R	
f CLEAR PRGM	000-
f LBL B	001-42,21,12
STO I	002- 44 25
R ↓	003- 33

* Program "B" exploits a tricky property of the **Σ-** and **Σ+** keys whereby certain calculations can be carried out to 13 significant digits before being rounded back to 10.

Keystrokes	Display
STO 0	004- 44 0
STO 8	005- 44 8
x\leqy	006- 34
STO 1	007- 44 1
STO 9	008- 44 9
f SCI 2	009-42, 8, 2
f LBL .8	010-42,21, .8
f CLEAR Σ	011- 42 32
RCL 8	012- 45 8
STO 7	013- 44 7
RCL \div I	014-45,10,25
g RND	015- 43 34
RCL I	016- 45 25
g Σ -	017- 43 49
RCL 9	018- 45 9
f x\leq 7	019-42, 4, 7
x\leqy	020- 34
RCL 8	021- 45 8
g Σ -	022- 43 49
R \downarrow	023- 33
g Σ -	024- 43 49
RCL 7	025- 45 7
g ABS	026- 43 16
RCL 9	027- 45 9
g ABS	028- 43 16
g x\leqy	029- 43 10
GTO B	030- 22 12
ENTER	031- 36
g R \uparrow	032- 43 33
STO 8	033- 44 8
RCL 7	034- 45 7
STO 9	035- 44 9
g ABS	036- 43 16
EEX	037- 26
2	038- 2
0	039- 0
x	040- 20
RCL 1	041- 45 1
g ABS	042- 43 16
g x\leqy	043- 43 10

Keystrokes	Display
GTO .8	044- 22 .8
f LBL B	045-42,21,12
f FIX 9	046-42, 7, 9
RCL 8	047- 45 8
g x²	048- 43 11
STO 7	049- 44 7
RCL I	050- 45 25
RCL 9	051- 45 9
g Σ-	052- 43 49
RCL 7	053- 45 7
g TEST 2	054-43,30, 2
GTO .7	055- 22 .7
√x	056- 11
f x_≧ 0	057-42, 4, 0
g TEST 2	058-43,30, 2
RCL - 0	059-45,30, 0
g TEST 3	060-43,30, 3
RCL + 0	061-45,40, 0
f x_≧ 1	062-42, 4, 1
g TEST 0	063-43,30, 0
RCL ÷ 1	064-45,10, 1
RCL 1	065- 45 1
RCL ÷ I	066-45,10,25
g RTN	067- 43 32
f LBL .7	068-42,21, .7
CHS	069- 16
√x	070- 11
RCL ÷ I	071-45,10,25
ENTER	072- 36
CHS	073- 16
RCL 0	074- 45 0
RCL I	075- 45 25
÷	076- 10
x_≧y	077- 34
f I	078- 42 25
ENTER	079- 36
g R↑	080- 43 33
f I	081- 42 25
g RTN	082- 43 32
g P/R	

This program's accuracy is phenomenal: better than nine significant digits even for the imaginary parts of nearly indistinguishable complex roots (as when $c = 4,877,163,849$ and $b = 4,877,262,613$ and $a = 4,877,361,379$); if the roots are integers, real or complex, and if $a = 1$, then the roots are calculated exactly (as when $c = 1,219,332,937 \times 10^1$, $b = 111,111.5$, and $a = 1$). But the program is costly; it uses more than twice as much memory for both program and data as does subroutine "A", and much more time, to achieve nine significant digits of accuracy instead of five in a few cases that can hardly ever matter—simply because the quadratic's coefficients can hardly ever be calculated exactly. If any coefficient c , b , or a is uncertain by as much as one unit in its 10th significant digit, then subroutine "B" is overkill. Subroutine "B" is like Grandmother's expensive chinaware, reserved for special occasions, leaving subroutine "A" for everyday use.

Index

Page numbers in **bold type** indicate primary references; page numbers in regular type indicate secondary references.

A

Absolute error, **173, 182**

Accuracy

in Complex mode, **73-75**

of integrand, **47-49**

of numerical calculations, **172-211**

of solutions to linear system, **103-104**

Aliasing, **46**

Analysis, discounted cash flow, **39-44**

Analysis of variance, **133-140**

Angle in triangle, **194-199**

Annuities, **26-39**

Annuity, ordinary, **27**

Annuity due, **27-28**

Annunciator, **C, 205**

Annunciator, trig mode, **68**

ANOVA table, **133, 134, 140**

Augmented matrix, **141**

Augmented normal equations, **111**

Augmented system, **142**

B

Backward error analysis, **187-211**

Balloon payment, **27, 29, 36**

Binomial theorem, **176**

Bounding search, **161, 162**

Branch, principal, **69-72**

Bridge too short, **174**

Broken calculator, **172, 175-176**

C

Calculation time, \boxed{f} , **49-55**

Calculations, numerical accuracy, **172-211**

Cancellation, **176-178, 200, 207**

Cash flow analysis, discounted, **39-44**
 Cash flow diagram, **28, 28-44**
 Characteristic equation, **148**
 Column norm, **99**
 Complementary error function, **60-64**
 Complementary normal distribution function, **60-64**
 Complex components, accurate, **74**
 Complex equations, solving large system, **128-131**
 Complex math functions, **68-72**
 Complex mode, **65-95**
 accuracy, **73-75**
 [SOLVE] and $\left[\frac{f}{f} \right]$, **73**
 Complex multivalued functions, **69-72**
 Complex number, n th roots, **69, 78-80**
 Complex number, storing and recalling, **76-78**
 Complex potential function, **89-95**
 Complex relative error, **183**
 Complex roots of equation, **16-17, 80-85**
 Complex roots of quadratic equation, **205-211**
 Complex single-valued functions, **69**
 Components, accurate complex, **74**
 Compound amounts, **26-39**
 Condition number, **98-102, 107, 201**
 Conformal mapping, **89**
 Constrained least-squares, **111, 115-116, 143**
 Consumer price index, **137-140, 147-148**
 Contour integral, **85-89**
 Correctly rounded result, **179-183**
 perturbed input, **184-211**
 Covariance matrix, **131**
 Critical point, **160, 162, 163**

D

Declination, **11-15**
 Decomposition, LU , **96-98, 117, 118**
 descriptor, **97**
 Deflation, **10**
 Degrees of freedom, **132**
 Delay equation, **81-85**
 Derivative, **10, 17-20, 192**
 Descartes' Rule of Signs, **10-11**

Descriptor of LU decomposition, **97**
 Determinant, **97-98, 118**
 Diagram, cash flow, **28, 28-44**
 Discounted cash flow analysis, **39-44**
 Discounted rate of return, **39**
 Display format, **45-46, 48**
 Doolittle method, **97**

E

Eigenvalue, **148-160**
 storage, **159-160**
 Eigenvector, **149, 154-160**
 Electrostatic field, **59**
 Endpoint, $\left[\frac{f}{z} \right]$ sampling at, **46-47, 56**
 Equations
 complex, solving large system, **128-131**
 equivalent, **9-10**
 solving inaccurate, **10**
 solving nonlinear system, **122-128**
 with several roots, **10**
 Equipotential line, **89-95**
 Equivalent equations, **9-10**
 Error 0, **29, 196, 199, 207**
 Error 1, **162, 167**
 Error 4, **29, 40**
 Error 8, **9, 23**
 Error analysis, backward, **187-211**
 Error function, **60-64**
 complementary, **60-64**
 Error, **173**
 absolute, **173, 182**
 hierarchy, **178**
 in matrix elements, **100-101**
 misconceptions, **172-178**
 relative, **173, 182, 183**
 Example
 angle in triangle, **194-199**
 annuities, **34-39**
 bridge too short, **174**
 broken calculator, **172, 175-176**
 cash flow, **43-44**
 compound amounts, **34-39**

consumer price index regression, 137-140, 147-148
 contour integral, 88-89
 declination of sun, 11-15
 delay equation, 81-85
 eigenvectors, 157-159
 equipotential line, 95
 field intensity of antenna, 17-25
 filter network, 128-131
 Gamma function, 65-68
 lunar phases, 186
 normal distribution function, 64
 n th roots of complex number, 80
 optimizing box, 168-171
 pennies, 173, 180-183
 pi, 173, 184-186
 quadratic surface, 153-154
 residual correction, 121
 roots of quadratic equation, 191, 205-211
 special functions, 64
 storing and recalling complex numbers, 77-78
 streamline, 93-94
 subdividing interval of integration, 51-54
 transformation of variables, 54-55
 unbiased test of hypothesis, 122-128
 Extended precision, 47, 104, 208
 Extremes of function, 17-25

F

F ratio, 132-140
 Factorization, orthogonal, 113-116, 140-148
 Field intensity, 17-25
 Financial equation, 29, 39
 Financial problems, 26-44
 Format, display, 45-46, 48
 Frobenius norm, 99
 Functions, complex, 68-73
 Future value, 26-39

G

Gamma function, complex, 65-68
 Gradient, 160, 162
 Grandmother's expensive chinaware, 211

H

Hierarchy of error, **178**
Horner's method, **11, 12**
Hyperbolic cylinder, **153-154**

I

Identity matrix, **119**
Ill-conditioned matrix, **98-102, 107, 155**
Ill-conditioned system of equations, **104-110**
Improper integral, **55-60**
Inaccurate equations, solving, **10**
Inaccurate roots, **9-10**
Input noise, **187-192**
Integral
 contour, **85-89**
 evaluating difficult, **55-60**
 improper, **55-60**
Integration, numerical, using $\boxed{\beta}$, **45-64**
Integration in Complex mode, **73**
Interchange, row, **97, 117**
Interest rate, **26-44**
Internal rate of return, **39-44**
Interval of integration, subdividing, **50-54, 58**
Interval reduction, **161, 162**
Inverse iteration, **155**
Inverse of function, **69**
Inverse of matrix, **98, 101-102, 110, 118, 187**
 backward error analysis, **200-204**
IRR, **39-44**
Iterative refinement, **103-104, 119-121**

J

Jordan canonical form, **155**

L

Large system of complex equations, solving, **128-131**
Least-squares, **110-116, 131-148, 187**
 linearly constrained, **111, 115-116, 143**
 weighted, **111, 115, 143**
Level 0, **178**
Level 1, **179-183, 190, 194**
Level 1C, **183**
Level 2, **184-211**

Level ∞ , 179
 Line search, 161
 Linear model, 131
 Linear regression, multiple, 131. *See also* Least-squares
 Linear system, accuracy of numerical solution, 103-104
 Linearly constrained least-squares, 111, 115-116, 143
 Lower-triangular matrix, 96
 LU decomposition, 96-98, 117, 118
 descriptor, 97
 Lunar phases, 186

M

Mapping, contour, 89
 Mathematical functions, complex, 68-72
 Mathematical functions, pure, 47-49
 Mathematical model, 48
 Matrix elements, errors in, 100-101
 Matrix inversion, backward error analysis, 200-204
 Matrix operations, 76-78, 96-171
 error levels, 178, 179, 192
 Maximum of function, 17-25, 160
 Mean-adjusted regression sum of squares, 134
 Minimum of function, 17-25, 160
 Model, linear, 131
 Model, mathematical, 48
 Monotonicity, 180, 186
 Multiple linear regression, 131. *See also* Least-squares
 Multiple root, 10
 Multivalued functions, complex, 69-72

N

Nearly singular matrix, 107, 117-118, 201, 204
 Net present value, 39-44
 equation, 39
 Network, filter, 128-131
 Newton's iteration method, 80-82, 122
 Noise, input and output, 187-192
 Nonlinear equations, solving system, 122-128
 Nonsingular matrix, 101-102, 117
 Norm, 99, 106, 200
 Normal distribution, 122-123, 132
 Normal distribution function, 48, 60-64
 complementary, 60-64

Normal equations, **110-113, 131-140**
 augmented, **111**
 weighted, **111**
NPV, **39-44**
 equation, **39**
*n*th roots of complex number, **69, 78-80**
Number of correct digits, **103, 121**
Numerical calculations, accuracy, **172-211**
Numerical integration, **45-64**
Numerical solutions to linear system, accuracy, **103-104**
Numerically finding roots, **6, 6-44**

O

Optimization, **160-171**
Ordinary annuity, **27**
Orthogonal factorization, **113-116, 140-148**
Orthogonal matrix, **113, 141, 142, 149**
Output noise, **188-192**
Overflow, **179**

P

Payment, **26-39**
Pennies, **173, 180-183**
Phases, lunar, **186**
Physical situations, **47-49**
Pi, **173, 184-186**
Pivots, **118**
Polar form, **68**
Polynomials, **10-15**
Potential function, complex, **89-95**
Precision, extended, **47, 104, 208**
Preconditioning a system, **107-110**
Present value, **26-44**
Principal branch, **69-72**
Principal value, **69-72**

Q

Quadratic equation, roots, **191, 205-211**
Quadratic surface, **149, 153-154**

R

Radians, used in Complex mode, **68**
Rate of return, **39-44**

Recalling complex numbers, **76-78**
 Rectangular form, **68**
 Refinement, iterative, **103-104, 119-121**
 Regression, multiple linear, **131**. *See also* Least-squares
 Regression sum of squares, **132-140**
 mean-adjusted, **134**
 Relative error, **173, 182, 183**
 complex, **73-75**
 Relative uncertainty of matrix, **100**
 Repeated estimation, **23-25**
 Residual, **103-104, 110, 132, 201**
 Residual correction, **103-104, 119-121**
 Residual sum of squares, **132-140**
 Resonance, **46**
 Return, rate of, **39-44**
 Romberg method, **46**
 Roots
 complex, **16-17**
 equations with several, **10**
 inaccurate, **9-10**
 multiple, **10**
 not found, **9, 29, 92**
 numerically finding, **6, 6-44**
 of complex number, **69, 78-80**
 of equation, complex, **80-85**
 of quadratic equation, **191, 205-211**
 Round-off error, **47, 49**. *See also* Rounding error
 Rounding error, **111, 113, 118, 172-211**
 Row interchange, **97, 117**
 Row norm, **99, 200**

S

Saddle-point, **162**
 Samples, \boxed{f} , **46-47, 50, 56, 73**
 Samples, $\boxed{\text{SOLVE}}$, **7-9, 73**
 Scaling a matrix, **104-107, 204**
 Scaling a system, **107**
 Secant method, **7**
 Sign change, **8**
 Sign symmetry, **180, 185**
 Single-valued functions, complex, **69**
 Singular matrix, **101-102, 117-118, 201**

Singularity and backward error analysis, 192-194
 Skew-symmetric matrix, 149
 Slope, 20-22
 Smaller root of quadratic equation, 191, 205-211
 Solutions to linear system, accuracy, 103-104
SOLVE, 6-44
 algorithm, 6-9, 73
 in Complex mode, 73
 Solving a system of equations, 15-17, 98, 100-101, 118, 122-128
 Solving a system of nonlinear equations, 122-128
 Solving equation for complex roots, 80-85
 Solving large system of complex equations, 128-131
 Steepest descent, 160
 Storing complex numbers, 76-78
 Streamline, 89-94
 Subdividing interval of integration, 50-54, 58
 Subinterval, 50-54
 Successive rows, 140-148
 Sum of squares, 132, 140
 Symmetric matrix, 148-149
 System of complex equations, solving large, 128-131
 System of equations, ill-conditioned, 104-110
 System of equations, solving, 15-17, 98, 100-101, 122-128
 System of nonlinear equations, solving, 122-128

T

Tail of function, 57-58
 Taylor series, 182
 Total sum of squares, 132-140
 Transformation of variables, 54-55
 Triangle, angle in, 194-199
 Trigonometric functions, 184-186
 Trigonometric modes, 68

U

Unbiased test, 122-123
 Uncertainty for β , 45-46
 Uncertainty of matrix, 100
 Unconstrained least-squares. *See* Least-squares
 Underflow, 50-51, 118, 179
 Upper-triangular matrix, 96, 113-114, 141

V

Variables, transforming, 54-55**W**

Weighted least-squares, 111, 115, 143**Weighted normal equations, 111****Y**

Yield, 39**Z**

Zero of polynomial, 10



Corvallis Division
1000 N.E. Circle Blvd., Corvallis, OR 97330, U.S.A.

00015-90011

Printed in U.S.A. 8/82