

# Topics Covered in Week 1

Instructor: Ms. Amber Shaikh

Course: CS- 2008 Numerical Computing

## Error Definition:

Numerical errors arise from the use of approximations to represent exact mathematical operations and quantities. For such errors, the relationship between the exact, or true, result and the approximation can be formulated as

$$\text{True value} = \text{approximation} + \text{error} \quad (4.1)$$

By rearranging Eq. (4.1), we find that the numerical error is equal to the discrepancy between the truth and the approximation, as in

$$E_t = \text{true value} - \text{approximation} \quad (4.2)$$

## Methods of Measuring Approximation Errors:

The following definition describes two methods for measuring approximation errors.

Suppose that  $p^*$  is an approximation to  $p$ . The **absolute error** is  $|p - p^*|$ , and the **relative error** is  $\frac{|p - p^*|}{|p|}$ , provided that  $p \neq 0$ .

To get the percent relative error just multiply the relative error by 100.

## Example no 1:

Determine the absolute and relative errors when approximating  $p$  by  $p^*$  when

- (a)  $p = 0.3000 \times 10^1$  and  $p^* = 0.3100 \times 10^1$ ;
  - (b)  $p = 0.3000 \times 10^{-3}$  and  $p^* = 0.3100 \times 10^{-3}$ ;
  - (c)  $p = 0.3000 \times 10^4$  and  $p^* = 0.3100 \times 10^4$ .
- 
- (a) For  $p = 0.3000 \times 10^1$  and  $p^* = 0.3100 \times 10^1$  the absolute error is 0.1, and the relative error is  $0.333\bar{3} \times 10^{-1}$ .
  - (b) For  $p = 0.3000 \times 10^{-3}$  and  $p^* = 0.3100 \times 10^{-3}$  the absolute error is  $0.1 \times 10^{-4}$ , and the relative error is  $0.333\bar{3} \times 10^{-1}$ .
  - (c) For  $p = 0.3000 \times 10^4$  and  $p^* = 0.3100 \times 10^4$ , the absolute error is  $0.1 \times 10^3$ , and the relative error is again  $0.333\bar{3} \times 10^{-1}$ .

This example shows that the same relative error,  $0.333\bar{3} \times 10^{-1}$ , occurs for widely varying absolute errors. As a measure of accuracy, the absolute error can be misleading and the relative error more meaningful, because the relative error takes into consideration the size of the value.

## Roundoff and Chopping Error:

The error that is produced when a calculator or computer is used to perform real-number calculations is called **round-off error**. It occurs because the arithmetic performed in a machine involves numbers with only a finite number of digits, with the result that calculations are performed with only approximate representations of the actual numbers. In a computer, only a relatively small subset of the real number system is used for the representation of all the real numbers. This subset contains only rational numbers, both positive and negative, and stores the fractional part, together with an exponential part.

## Example no 2:

Determine the five-digit (a) chopping and (b) rounding values of the irrational number  $\pi$ .

The number  $\pi$  has an infinite decimal expansion of the form  $\pi = 3.14159265 \dots$ . Written in normalized decimal form, we have

$$\pi = 0.314159265 \dots \times 10^1.$$

- (a) The floating-point form of  $\pi$  using five-digit chopping is

$$fl(\pi) = 0.31415 \times 10^1 = 3.1415.$$

- (b) The sixth digit of the decimal expansion of  $\pi$  is a 9, so the floating-point form of  $\pi$  using five-digit rounding is

$$fl(\pi) = (0.31415 + 0.00001) \times 10^1 = 3.1416.$$

## Computer Representation of numbers:

In 1985, the IEEE (Institute for Electrical and Electronic Engineers) published a report called *Binary Floating Point Arithmetic Standard 754–1985*. An updated version was published in 2008 as *IEEE 754-2008*. This provides standards for binary and decimal floating point numbers, formats for data interchange, algorithms for rounding arithmetic operations, and for the handling of exceptions. Formats are specified for single, double, and extended precisions, and these standards are generally followed by all microcomputer manufacturers using floating-point hardware.

A 64-bit (binary digit) representation is used for a real number. The first bit is a sign indicator, denoted  $s$ . This is followed by an 11-bit exponent,  $c$ , called the **characteristic**, and a 52-bit binary fraction,  $f$ , called the **mantissa**. The base for the exponent is 2.

Watch the given video for a 32-bit(binary digit) representation of a real number in IEEE 754.

<https://youtu.be/8afbTaA-gQQ>

# Underflow and Overflow error:

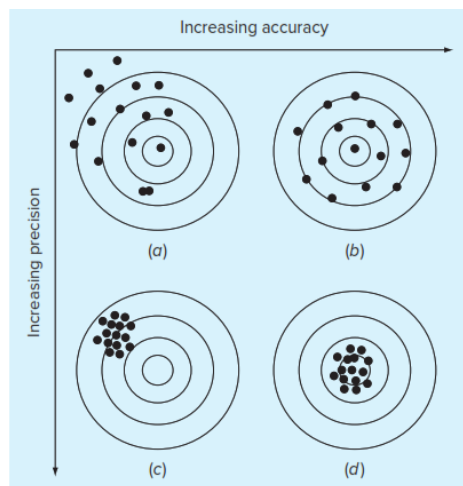
From the definition of floating-point number, there are **upper** and **lower** limits for the magnitudes of the numbers that can be expressed in a floating-point form.

Attempts to create numbers

- that are too small  $\Rightarrow$  **underflow** errors: the default option is to set the number to zero and proceed
- that are too large  $\Rightarrow$  **overflow** errors: generally fatal errors on most computers. With the IEEE floating-point format, overflow errors can be carried along as having a value of  $\pm\infty$  or *NaN*, depending on the context. Usually, an overflow error is an indication of a more significant problem or error in the program.

## Accuracy and Precision:

The errors associated with both calculations and measurements can be characterized with regard to their accuracy and precision. *Accuracy* refers to how closely a computed or measured value agrees with the true value. *Precision* refers to how closely individual computed or measured values agree with each other.



**FIGURE 4.1**

An example from marksmanship illustrating the concepts of accuracy and precision: (a) inaccurate and imprecise, (b) accurate and imprecise, (c) inaccurate and precise, and (d) accurate and precise.

# Significant Digits and Loss of Significance:

## RULES FOR SIGNIFICANT FIGURES

1. **All non-zero numbers ARE significant.** The number 33.2 has THREE significant figures because all of the digits present are non-zero.

2. **Zeros between two non-zero digits ARE significant.** 2051 has FOUR significant figures. The zero is between a 2 and a 5.

3. **Leading zeros are NOT significant.** They're nothing more than "place holders." The number 0.54 has only TWO significant figures. 0.0032 also has TWO significant figures. All of the zeros are leading.

4. **Trailing zeros to the right of the decimal ARE significant.** There are FOUR significant figures in 92.00.

92.00 is different from 92: a scientist who measures 92.00 milliliters knows his value to the nearest 1/100th milliliter; meanwhile his colleague who measured 92 milliliters only knows his value to the nearest 1 milliliter. It's important to understand that "zero" does not mean "nothing." Zero denotes actual information, just like any other number. You cannot tag on zeros that aren't certain to belong there.

5. **Trailing zeros in a whole number with the decimal shown ARE significant.** Placing a decimal at the end of a number is usually not done. By convention, however, this decimal indicates a significant zero. For example, "540." indicates that the trailing zero IS significant; there are THREE significant figures in this value.

6. **Trailing zeros in a whole number with no decimal shown are NOT significant.** Writing just "540" indicates that the zero is NOT significant, and there are only TWO significant figures in this value.

The number  $p^*$  is said to approximate  $p$  to  $t$  **significant digits** (or figures) if  $t$  is the largest nonnegative integer for which

$$\frac{|p - p^*|}{|p|} \leq 5 \times 10^{-t}.$$

## Example 3 and 4:

Suppose that  $x = \frac{5}{7}$  and  $y = \frac{1}{3}$ . Use five-digit chopping for calculating  $x + y$ ,  $x - y$ ,  $x \times y$ , and  $x \div y$ .

Note that

$$x = \frac{5}{7} = 0.\overline{714285} \quad \text{and} \quad y = \frac{1}{3} = 0.\overline{3}$$

implies that the five-digit chopping values of  $x$  and  $y$  are

$$fl(x) = 0.71428 \times 10^0 \quad \text{and} \quad fl(y) = 0.33333 \times 10^0.$$

Thus

$$\begin{aligned} x \oplus y &= fl(fl(x) + fl(y)) = fl(0.71428 \times 10^0 + 0.33333 \times 10^0) \\ &= fl(1.04761 \times 10^0) = 0.10476 \times 10^1. \end{aligned}$$

The true value is  $x + y = \frac{5}{7} + \frac{1}{3} = \frac{22}{21}$ , so we have

$$\text{Absolute Error} = \left| \frac{22}{21} - 0.10476 \times 10^1 \right| = 0.190 \times 10^{-4}$$

and

$$\text{Relative Error} = \left| \frac{0.190 \times 10^{-4}}{22/21} \right| = 0.182 \times 10^{-4}.$$

$$| \quad \lll / \lll \quad |$$

Table 1.2 lists the values of this and the other calculations.

Operation	Result	Actual value	Absolute error	Relative error
$x \oplus y$	$0.10476 \times 10^1$	$22/21$	$0.190 \times 10^{-4}$	$0.182 \times 10^{-4}$
$x \ominus y$	$0.38095 \times 10^0$	$8/21$	$0.238 \times 10^{-5}$	$0.625 \times 10^{-5}$
$x \otimes y$	$0.23809 \times 10^0$	$5/21$	$0.524 \times 10^{-5}$	$0.220 \times 10^{-4}$
$x \oslash y$	$0.21428 \times 10^1$	$15/7$	$0.571 \times 10^{-4}$	$0.267 \times 10^{-4}$

The maximum relative error for the operations in Example 3 is  $0.267 \times 10^{-4}$ , so the arithmetic produces satisfactory five-digit results. This is not the case in the following example.



Suppose that in addition to  $x = \frac{5}{7}$  and  $y = \frac{1}{3}$  we have

$$u = 0.714251, \quad v = 98765.9, \quad \text{and} \quad w = 0.111111 \times 10^{-4},$$

so that

$$fl(u) = 0.71425 \times 10^0, \quad fl(v) = 0.98765 \times 10^5, \quad \text{and} \quad fl(w) = 0.11111 \times 10^{-4}.$$

Determine the five-digit chopping values of  $x \ominus u$ ,  $(x \ominus u) \oplus w$ ,  $(x \ominus u) \otimes v$ , and  $u \oplus v$ .

These numbers were chosen to illustrate some problems that can arise with finite-digit arithmetic. Because  $x$  and  $u$  are nearly the same, their difference is small. The absolute error for  $x \ominus u$  is

$$\begin{aligned} |(x - u) - (x \ominus u)| &= |(x - u) - (fl(fl(x) - fl(u)))| \\ &= \left| \left( \frac{5}{7} - 0.714251 \right) - (fl(0.71428 \times 10^0 - 0.71425 \times 10^0)) \right| \\ &= |0.347143 \times 10^{-4} - fl(0.00003 \times 10^0)| = 0.47143 \times 10^{-5}. \end{aligned}$$

This approximation has a small absolute error, but a large relative error

$$\left| \frac{0.47143 \times 10^{-5}}{0.347143 \times 10^{-4}} \right| \leq 0.136.$$

The subsequent division by the small number  $w$  or multiplication by the large number  $v$  magnifies the absolute error without modifying the relative error. The addition of the large and small numbers  $u$  and  $v$  produces large absolute error but not large relative error. These calculations are shown in Table 1.3.

Operation	Result	Actual value	Absolute error	Relative error
$x \ominus u$	$0.30000 \times 10^{-4}$	$0.34714 \times 10^{-4}$	$0.471 \times 10^{-5}$	0.136
$(x \ominus u) \oplus w$	$0.27000 \times 10^1$	$0.31242 \times 10^1$	0.424	0.136
$(x \ominus u) \otimes v$	$0.29629 \times 10^1$	$0.34285 \times 10^1$	0.465	0.136
$u \oplus v$	$0.98765 \times 10^5$	$0.98766 \times 10^5$	$0.161 \times 10^1$	$0.163 \times 10^{-4}$

One of the most common error-producing calculations involves the cancelation of significant digits due to the subtraction of nearly equal numbers.

# Truncation Error:

## Taylor Polynomial:

Suppose  $f \in C^n[a, b]$ , that  $f^{(n+1)}$  exists on  $[a, b]$ , and  $x_0 \in [a, b]$ . For every  $x \in [a, b]$ , there exists a number  $\xi(x)$  between  $x_0$  and  $x$  with

$$f(x) = P_n(x) + R_n(x),$$

where

$$\begin{aligned} P_n(x) &= f(x_0) + f'(x_0)(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2 + \cdots + \frac{f^{(n)}(x_0)}{n!}(x - x_0)^n \\ &= \sum_{k=0}^n \frac{f^{(k)}(x_0)}{k!}(x - x_0)^k \end{aligned}$$

and

$$R_n(x) = \frac{f^{(n+1)}(\xi(x))}{(n+1)!}(x - x_0)^{n+1}.$$

Here  $P_n(x)$  is called the  **$n$ th Taylor polynomial** for  $f$  about  $x_0$ , and  $R_n(x)$  is called the **remainder term** (or **truncation error**) associated with  $P_n(x)$ . Since the number  $\xi(x)$  in the truncation error  $R_n(x)$  depends on the value of  $x$  at which the polynomial  $P_n(x)$  is being evaluated, it is a function of the variable  $x$ . However, we should not expect to be able to explicitly determine the function  $\xi(x)$ . Taylor's Theorem simply ensures that such a function exists, and that its value lies between  $x$  and  $x_0$ . In fact, one of the common problems in numerical methods is to try to determine a realistic bound for the value of  $f^{(n+1)}(\xi(x))$  when  $x$  is in some specified interval.

The infinite series obtained by taking the limit of  $P_n(x)$  as  $n \rightarrow \infty$  is called the **Taylor series** for  $f$  about  $x_0$ . In the case  $x_0 = 0$ , the Taylor polynomial is often called a **Maclaurin polynomial**, and the Taylor series is often called a **Maclaurin series**.



$$\begin{aligned}
 e^x &= 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \frac{x^4}{4!} + \dots \\
 \sin x &= x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \dots \\
 \cos x &= 1 - \frac{x^2}{2!} + \frac{x^4}{4!} - \frac{x^6}{6!} + \dots \\
 \ln(1+x) &= x - \frac{x^2}{2} + \frac{x^3}{3} - \frac{x^4}{4} + \dots \quad \text{for } |x| < 1 \\
 \tan^{-1}(x) &= x - \frac{x^3}{3} + \frac{x^5}{5} - \frac{x^7}{7} + \dots \quad \text{for } |x| < 1
 \end{aligned}$$

**Problem Statement.** In mathematics, functions can often be represented by infinite series. For example, the exponential function can be computed using

$$e^x = 1 + x + \frac{x^2}{2} + \frac{x^3}{3!} + \dots + \frac{x^n}{n!} \quad (\text{E4.1.1})$$

Thus, as more terms are added in sequence, the approximation becomes a better and better estimate of the true value of  $e^x$ . Equation (E4.1.1) is called a *Maclaurin series expansion*.

Starting with the simplest version,  $e^x = 1$ , add terms one at a time in order to estimate  $e^{0.5}$ . After each new term is added, compute the true and approximate percent relative errors with Eqs. (4.3) and (4.5), respectively. Note that the true value is  $e^{0.5} = 1.648721 \dots$ . Add terms until the absolute value of the approximate error estimate  $\epsilon_a$  falls below a prespecified error criterion  $\epsilon_s$  conforming to three significant figures.

**Solution.** First, Eq. (4.7) can be employed to determine the error criterion that ensures a result that is correct to at least three significant figures:

$$\epsilon_s = (0.5 \times 10^{2-3})\% = 0.05\%$$

Thus, we will add terms to the series until  $\varepsilon_a$  falls below this level.

The first estimate is simply equal to Eq. (E4.1.1) with a single term. Thus, the first estimate is equal to 1. The second estimate is then generated by adding the second term as in

$$e^x = 1 + x$$

or for  $x = 0.5$

$$e^{0.5} = 1 + 0.5 = 1.5$$

This represents a true percent relative error of [Eq. (4.3)]

$$\varepsilon_t = \left| \frac{1.648721 - 1.5}{1.648721} \right| \times 100\% = 9.02\%$$

Equation (4.5) can be used to determine an approximate estimate of the error, as in

$$\varepsilon_a = \left| \frac{1.5 - 1}{1.5} \right| \times 100\% = 33.3\%$$

Because  $\varepsilon_a$  is not less than the required value of  $\varepsilon_s$ , we would continue the computation by adding another term,  $x^2/2!$ , and repeating the error calculations. The process is continued until  $|\varepsilon_a| < \varepsilon_s$ . The entire computation can be summarized as

Terms	Result	$\varepsilon_t, \%$	$\varepsilon_a, \%$
1	1	39.3	
2	1.5	9.02	33.3
3	1.625	1.44	7.69
4	1.645833333	0.175	1.27
5	1.648437500	0.0172	0.158
6	1.648697917	0.00142	0.0158

Thus, after six terms are included, the approximate error falls below  $\varepsilon_s = 0.05\%$ , and the computation is terminated. However, notice that, rather than three significant figures, the result is accurate to five! This is because, for this case, both Eqs. (4.5) and (4.7) are conservative. That is, they ensure that the result is at least as good as they specify. Although, this is not always the case for Eq. (4.5), it is true most of the time.

**For graphical visualization of series plot function kindly see the attached pdf.**

