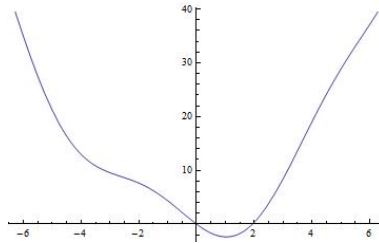


INTRODUCTION

1 Introduction

- Numerical methods are mathematical techniques used for solving mathematical problems that cannot be solved or are difficult to solve (example: eq.1). The numerical solution is an approximate numerical value for the solution. Although numerical solutions are an approximation, they can be very accurate.
 - Example: Find the roots of the following equation

$$f(x) = x^2 - 4 \sin(x) = 0$$



- In many numerical methods, the calculations are executed in an iterative manner until a desired accuracy is achieved.
 - Example: start at one value of x then change its value in small increment. A change in the sign of $f(x)$ indicates that there is a root within the last increment.
- Today, numerical methods are used in fast electronic digital computers that make it possible to execute many tedious and repetitive calculations that produce accurate (even though not exact) solutions in a very short time.
- For every type of mathematical problem there are several numerical techniques that can be used. The techniques differ in accuracy, length of calculations, and difficulty in programming.
 - Example: to solve a nonlinear equation of the form $f(x) = 0$ one can use among others
 - * Bisection method
 - * Regula Falsi method
 - * Newton's method
 - * Secant method

These techniques and others will be discussed in the next chapters.

2 Errors in numerical solutions

Since numerical solutions are an approximation, and since the computer program that executes the numerical method might have errors, a numerical solution needs to be examined closely. There are three major sources of error in computation: **human errors**, **truncation errors**, and **round-off errors**.

2.1 Human errors

Typical human errors are arithmetic errors, and/or programming errors: These errors can be very hard to detect unless they give obviously incorrect solution. In discussing errors, we shall assume that human errors are not present.

- Example of arithmetic errors: When parentheses or the rules about orders of operation are misunderstood or ignored:
 - * You can remember the correct order of operations rules by the mnemonic **PEMDAS**, which says to compute anything inside Parentheses first, then compute Exponential expressions (powers) next, then compute Multiplications and Divisions from left to right, and finally compute Additions and Subtractions from left to right. The highest priority for parentheses means that you should follow the remaining rules for anything inside the parentheses to arrive at a result for that part of the calculation.

2.2 Truncation and Round-off errors

- Error in computation is the difference between the exact answer X_{ex} and the computed answer X_{cp} . This is also known as true error

$$true\ error = X_{cp} - X_{ex}$$

- Since we are usually interested in the magnitude or absolute value of the error we define

$$true\ absolute\ error = |X_{cp} - X_{ex}|$$

- The absolute error can sometimes be misleading. Even if the absolute error is small in magnitude, the approximation may still be grossly inaccurate if the exact value X_{ex} is even smaller in magnitude. For this reason, it is preferable to measure accuracy in terms of the relative error.

$$true\ relative\ error = \frac{|X_{cp} - X_{ex}|}{|X_{ex}|}$$

- Note that the errors defined above cannot be determined in problems that require numerical methods for their solution, this is because the exact solution X_{ex} is not known. These error quantities are useful for evaluating the accuracy of different numerical methods when the exact solution is known (problem solved analytically).
- Since the true errors cannot, in most cases, be calculated, other means are used for estimating the accuracy of a numerical solution. This will be discussed in more details in later chapters

2.2.1 Truncation error:

Numerical methods use approximations for solving problems. The errors introduced by the approximations are the truncation errors. For example, consider the Taylor series expansion

$$e^x = 1 + x + \frac{x^2}{2!} + \dots + \frac{x^n}{n!} + \dots$$

If the formula is used to calculate $e^{0.3}$ we get

$$e^{0.3} = 1 + 0.3 + \frac{0.3^2}{2!} + \dots + \frac{0.3^n}{n!} + \dots$$

Where do we stop the calculation? How many terms do we include? Theoretically the calculation will never stop. There are always more terms to add on. If we do stop after a finite number of terms, we will not get the exact answer. For example if we do take the first four terms as the approximation we get

$$r = e^{0.3} \approx 1 + 0.3 + \frac{0.3^2}{2!} + \frac{0.3^3}{3!} = \bar{r}$$

For this calculation, the truncation error is

$$\text{truncation error} = r - \bar{r}$$

Another example of truncation error that is the approximate calculation of derivatives. The value of the derivative of a function at a point can be approximated by the expression:

$$\frac{df}{dx} \Big|_{x=x_1} = \frac{f(x_2) - f(x_1)}{x_2 - x_1}$$

where x_2 is a point near x_1 . The difference between the value of the true derivative and the value that is calculated with this equation is called a truncation error. The truncation error is dependent on the specific numerical method or algorithm used to solve a problem. The truncation error is independent of round-off error.

2.2.2 Round-off error

Numbers can be represented in various forms. The familiar decimal system (base 10) uses ten digits 0, 1, ..., 9. A number is written by a sequence of digits that correspond to multiples of powers of 10.

* Example:

$$\begin{array}{cccccccccc}
 10^4 & 10^3 & 10^2 & 10^1 & 10^0 & 10^{-1} & 10^{-2} & 10^{-3} & 10^{-4} & \\
 \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \\
 6 & 0 & 7 & 2 & 4 & . & 3 & 1 & 2 & 5 \\
 6 \times 10^4 + 0 \times 10^3 + 7 \times 10^2 + 2 \times 10^1 + 4 \times 10^0 + 3 \times 10^{-1} + 1 \times 10^{-2} + 2 \times 10^{-3} + 5 \times 10^{-4} = 60,724.3125
 \end{array}$$

In general, however, a number can be represented using other bases. A form that can be easily implemented in computers is the binary (base 2) system. In the binary system, a number is represented by using the two digits 0 and 1. A number is then written as a sequence of zeros and ones that correspond to multiples of powers of 2.

* Example:

$$\begin{array}{cccccccc}
 2^4 & 2^3 & 2^2 & 2^1 & 2^0 & 2^{-1} & 2^{-2} & 2^{-3} \\
 \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow \\
 1 & 0 & 0 & 1 & 1 & . & 1 & 0 & 1 \\
 1 \times 2^4 + 0 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 1 \times 2^0 + 1 \times 2^{-1} + 0 \times 2^{-2} + 1 \times 2^{-3} \\
 1 \times 16 + 0 \times 8 + 0 \times 4 + 1 \times 2 + 1 \times 1 + 1 \times 0.5 + 0 \times 0.25 + 1 \times 0.125 = 19.625
 \end{array}$$

Computers store and process numbers in binary (base 2) form. Each binary digit (one or zero) is called a bit (for binary digit). Binary arithmetic is used by computers because modern transistors can be used as extremely fast switches. Therefore, a network of these may be used to represent strings of numbers with the 1 referring to the switch being in the on position and 0 referring to the off position. Various operations are then performed on these sequences of ones and zeros.

Computer memory has a finite capacity. This obvious fact has far reaching implications on the representation of real numbers, which in general do not have a finite uniform representation. To accommodate large and small numbers, real numbers are written in floating point representation. A floating-point number is typically expressed in the scientific notation, with a fraction (F), and an exponent (e) of a certain radix (r), in the form of $F \times r^e$. Decimal numbers use radix of 10 ($F \times 10^e$); while binary numbers use radix of 2 ($F \times 2^e$).

- * Upon carefully examining the possible range of exponents you will notice that it is not fully utilized. For double precision $2^{11} = 2048$ different exponents could be distinguished, and for single precision $2^8 = 256$ are available. But only 2046 and 254, respectively, are used in practice. This is because the IEEE standard reserves the endpoints of the exponent range for special purposes. The largest number precisely representable in a standard long word is therefore.

$$\left[1 + \sum_{i=1}^{52} \frac{1}{2^i} \right] \times 2^{1023} \approx 1.8 \times 10^{308}$$

attempts to define a number with $|number|$ larger than this causes overflow error and the smallest positive number is

$$[1 + 0] \times 2^{-1022} \approx 2.2 \times 10^{-308}$$

attempts to define a number with $|number|$ smaller than this causes underflow error
How are 0 and 1 stored? Here is where the endpoints of the exponent are used, and the conventions are simple and straightforward:

- For 0, set all binaries in the exponent and mantissa to zero, with the sign s arbitrary; i.e., the minimal positive value representable in the system is considered 0.

$$[1 + 0] \times 2^{-1023} \approx 1.1 \times 10^{-308}$$

- For infinity, set the binaries in the exponent and mantissa to 1

$$\left[1 + \sum_{i=1}^{52} \frac{1}{2^i} \right] \times 2^{1024} \approx 1.8 \times 10^{308}$$

- The pattern in which all binaries in the exponent are 1 and one or more binaries in the mantissa is zero is by convention NaN.
- In single precision, or short word, the largest number precisely representable is

$$\left[1 + \sum_{i=1}^{23} \frac{1}{2^i} \right] \times 2^{127} \approx 3.4 \times 10^{38}$$

and the smallest positive number is

$$[1 + 0]2^{-126} \approx 1.210^{-38}$$

- * When overflow occurs in the course of a calculation, this is generally fatal. But underflow is non-fatal: the system usually sets the number to 0 and continues. (Matlab does this, quietly.)

