

Chapter 1: Computer Arithmetic

Contents

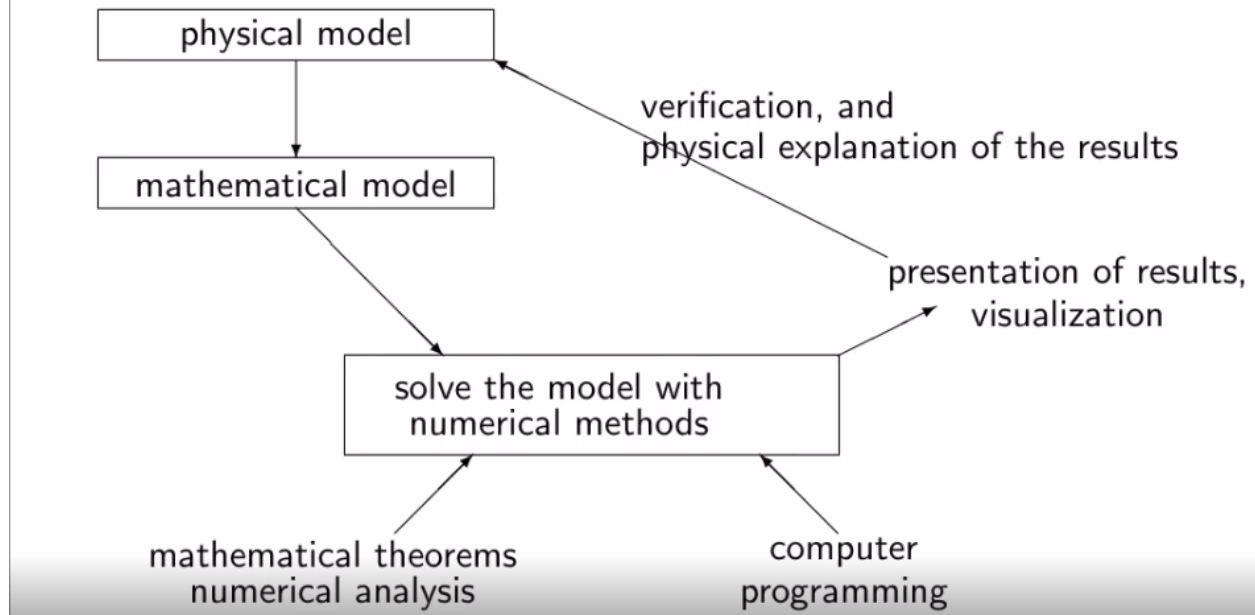
1	Introduction to Numerical Computation	1
1.1	Numerical Methods	1
2	Representation of numbers in different bases	2
2.1	Examples	3
3	Floating Point Representation	5
4	Loss of Significance	5
5	Review of Taylor Series	5
6	Finite Difference Approximation	5

1 Introduction to Numerical Computation

1.1 Numerical Methods

They are **algorithms** that compute **approximations** to functions, their derivatives, their integrations, and solutions to various equations etc. Such algorithms could be programmed on a computer.

Below is an overview on how various aspects are related.



Numerical methods are not about numbers. It is about mathematical ideas and insight. A little idea can go a long way. Some classical problems:

- Development of algorithms
- Implementation
- Some analysis, including error-estimates, convergence, stability, etc.

Matlab

2 Representation of numbers in different bases

There have been several different ways of representing numbers:

- 10: decimal, daily use
- 2: binary, computer use
- 8: octal
- 16: hexadecimal, ancient China
- 20: vigesimal, ancient France (can be seen in french language)
- 60: sexagesimal, used by Babylonians

In principle we can use any number β as the base. Writing such a number with a decimal point gives us an integer part and fraciton part.

$$\begin{array}{l}
\begin{array}{cc}
\text{integer part} & \text{fractional part} \\
\left(\overbrace{a_n a_{n-1} \cdots a_1 a_0} \cdot \overbrace{b_1 b_2 b_3 \cdots} \right)_{\beta} & \\
= & a_n \beta^n + a_{n-1} \beta^{n-1} + \cdots + a_1 \beta + a_0 \quad \text{(integer part)} \\
& + b_1 \beta^{-1} + b_2 \beta^{-2} + b_3 \beta^{-3} + \cdots \quad \text{(fractional part)}
\end{array}
\end{array}$$

This above formula allows us to convert a number in any base β into base 10.

2.1 Examples

2.1.1 Octal to Decimal

We have $(45.12)_8$.

$$4 * 8^1 + 5 * 8^0 + 1 * 8^{-1} + 2 * 8^{-2} = (37.15625)_{10}$$

2.1.2 Octal to Binary

$$\begin{array}{lcl} (1)_8 & = & (1)_2 \\ (2)_8 & = & (10)_2 \\ (3)_8 & = & (11)_2 \\ (4)_8 & = & (100)_2 \\ (5)_8 & = & (101)_2 \\ (6)_8 & = & (110)_2 \\ (7)_8 & = & (111)_2 \\ (10)_8 & = & (1000)_2 \end{array}$$

Since 8 is a power of 2 it makes it very easy to convert from octal to binary:

$$(5034)_8 = (101000011100)_2$$

Each digit in octal gets converted to 3 digits in binary. And vice versa...

You can also go back from binary to octal:

$$(110010111001)_2 = (6271)_8$$

2.1.3 Decimal to Binary

Write $(12.45)_{10}$ in binary. (base 2) This is particularly interesting because we use decimal and computer uses binary. This conversion takes two steps. First we convert the integer part into binary.

Procedure: Divide the integer by 2 and store the remainder of each step until integer is zero. (Euclid's algorithm.)

		<i>(remainder)</i>	
2	<u>12</u>	0	
2	<u>6</u>	0	
2	<u>3</u>	1	
2	<u>1</u>	1	
	0		

$$\Rightarrow (12)_{10} = (1100)_2$$

Now to convert the fractional part to binary you multiply by 2 and store the integer part for the result.

0.45	×	2	
<u>0.9</u>	×	2	
<u>1.8</u>	×	2	
<u>1.6</u>	×	2	
<u>1.2</u>	×	2	
<u>0.4</u>	×	2	
<u>0.8</u>	×	2	
<u>1.6</u>	×	2	
...			

$$\Rightarrow (0.45)_{10} = (0.01110011001100\dots)_2.$$

Note that the fractional part here is not finite. Putting them together:

$$(12.45)_{10} = (1100.01110011001100\dots)_2$$

So how do we store this kind of a number in a computer?

3 Floating Point Representation

4 Loss of Significance

5 Review of Taylor Series

6 Finite Difference Approximation