

A Study on some Properties of Floating-point Systems

By Myoung Lyoul Kim & Byung Chul Kim

Cheonbug National University, Cheonju, Korea

§1. Introduction

An agent which performs the algorithm of the given problem is called a *processor*. The processor may be a man or a machine. Let M be the set of machine representable numbers. Each number in set M can have only finite number of digits in it.

It is natural to postulate that the approximation of any number $x \notin M$ by a machine number $rd(x) \in M$ should satisfy $|x - rd(x)| \leq |x - g|$ for all $g \in M$.

Such a machine number approximation $rd(x)$ can be obtained in most case by *rounding*. [1] Generally, machine arithmetic number systems can be divided into the five categories: [2]

- ① Conventional Radix Number System ② Signed-Digit Number System
- ③ Residue Number System ④ Rational Number System ⑤ Logarithmic Number System

This paper shows the some properties of the number system on general-purpose digital computers by referencing [3] mainly.

However, the computer arithmetic is not dealt in this paper.

§2. Floating-point Systems

On computer memory, a given number may be stored in of two modes: *fixed-point and floating-point*. A floating-point hardware capabilities are now used for the vast majority of all numerical computation on computers because it has solved the problems in fixed-point hardware. Floating-point can represent numbers whose range is very large but it is handled with difficulty.

By normalizing the b -adic representation, we employ the symbol R_b to denote all real numbers, including zero:

$$R_b := \{0\} \cup \{x = Sm \cdot b^e \mid S \in \{+, -\}, b \in \mathbb{N}, b > 1, e \in \mathbb{Z}, m = \sum_{i=1}^{\infty} x[i]b^{-i}, \\ x[i] \in \{0, 1, \dots, b-1\}, x[1] \neq 0, x[i] \leq b-2 \text{ for infinitely many } i\}.$$

In general the element of R_b cannot be represented on computer. Only truncated version of these elements can be so represented.

Definition. A real number is called a *normalized floating-point number* if it is an element of the following sets $S_{b,l}$ or $S = S(b, l, e1, e2)$. These sets are called *floating-point systems*:

- (i) $S_{b,l} := \{x \in R_b \mid m = \sum_{i=1}^l x[i]b^{-i}\}.$
- (ii) $S = S(b, l, e1, e2) := \{x \in R_b \mid e1 \leq e \leq e2, e1, e, e2 \in \mathbb{Z}\}.$

In order to have a unique representation of zero available in S , we put additionally that $\text{sgn}(0) = +$, $\text{mant}(0) = 0.00 \dots 0$ (l zeros after the b -ary point), and $\text{exp}(0) = e1$.

§3. Properties.

The following properties are easily shown from the above definition.

Property 1 S is symmetric, i.e., $0, 1 \in S \wedge -x \in S$ for all $x \in S$.

Property 2 $S_{b,l}$ and S are bounded, and we have $S \subset S_{b,l} \subset R_b = R$.

Especially, $\sup S = B$ and $\inf S = -B$, where $B = +0$. $(b-1)(b-1)\dots(b-1) \cdot b^{\epsilon_1}$.

Property 3 (i) $S_{b,l}$ is countable, i.e., $\|S_{b,l}\| = \aleph_0$.

(ii) S is finite, and its element represents a rational number.

The representation is unique. $\|S\| = 2(b-1)b^{l-1}(e_2 - e_1 + 1) + 1$.

Property 4 The floating-point systems are not dense. The floating-point number in S are not uniformly distributed between $[-B, -L]$ and $[L, B]$, where $L = 0.10\dots 0b^{\epsilon_1}$. Of course the floating-point systems are not continuum. They are discrete.

To continue, we introduce the following notation:

$$R^* := R \cup \{-\infty\} \cup \{+\infty\}, \quad S^* := S \cup \{-\infty\} \cup \{+\infty\}.$$

Let \circ denote the rounding to the nearest floating-point numbers of S^* , i.e., $\circ: R^* \rightarrow S^*$. An we define the floating-point operations as: For all $x, y \in S^*$ $x * y := \circ(x * y)$, where $*$ in S^* is a operation corresponding to a usual operation $* \in \{+, -, \times, \div\}$ in R^* .

We see easily that S^* is not closed under its operation $*$.

Property 5 \circ is not homomorphic.

For example, $b := 10$, $m := -0.9(0.1)$ 0.9 , $e \in \{-1, 0, 1\}$ and $x := 0.34$, $y := 0.54 \in R^*$.

Then we get $\circ x = 0.3$, $\circ y = 0.5$; $(\circ x) \oplus (\circ y) = \circ(\circ x + \circ y) = \circ(0.8) = 0.8$, $\circ(x + y) = \circ(0.88) = 0.9$, i.e., $(\circ x) \oplus (\circ y) \neq \circ(x + y)$.

Property 6 S^* is not associative for \oplus and \otimes .

With $x := 0.7$, $y := 0.7$, $z := 0.9$, we obtain

$$(x \oplus y) \oplus z = \circ(1.4) \oplus 0.9 = 0.1 \cdot 10 \oplus 0.9 = \circ(0.19 \cdot 10) = 0.2 \cdot 10,$$

$$x \oplus (y \oplus z) = 0.7 \oplus (\circ 1.6) = 0.7 \oplus 0.2 \cdot 10 = \circ(0.27 \cdot 10) = 0.3 \cdot 10, \text{ i.e., } (x \oplus y) \oplus z \neq x \oplus (y \oplus z);$$

$$\text{and } (x \otimes y) \otimes z = (\circ 0.49) \otimes 0.9 = 0.5 \otimes 0.9 = \circ(0.45) = 0.5 \times 10^0,$$

$$x \otimes (y \otimes z) = 0.7 \otimes (\circ 0.63) = 0.7 \otimes 0.6 = \circ(0.42) = 0.4 \times 10^0, \text{ i.e., } (x \otimes y) \otimes z \neq x \otimes (y \otimes z).$$

Property 7 S^* is not distributive for \otimes over \oplus .

For $x := 0.3$, $y := 0.7$, $z := 0.9$, we obtain $x \otimes (y \oplus z) = 0.3 \otimes (\circ 1.6) = 0.3 \otimes 0.2 \cdot 10 = 0.6$,

$$x \otimes y \oplus x \otimes z = (\circ 0.21) \oplus (\circ 0.27) = 0.2 \oplus 0.3 = 0.5, \text{ i.e., } x \otimes (y \oplus z) \neq x \otimes y \oplus x \otimes z.$$

References

1. J. Stoer, R. Bulirsch: *Introduction to Numerical Analysis*; Springer-Verlag New York, Inc., 1980 p. 5.
2. Kai Hwang: *Computer Arithmetic (Principles, Architecture, and Design)*; John Wiley & Son: Inc., 1979, p. 4.
3. Ulrich W. Kulisch, Williard L. Miranker: *Computer Arithmetic in Theory and Practice*; Academic Press, Inc., 1981, pp. 149-157.