

# B

## FLOATING-POINT NUMBERS

In many calculations the range of numbers used is very large. For example, a calculation in astronomy might involve the mass of the electron,  $9 \times 10^{-28}$  grams, and the mass of the sun,  $2 \times 10^{33}$  grams, a range exceeding  $10^{60}$ . These numbers could be represented by

```
00000000000000000000000000000000.000000000000000000000000000009  
2000000000000000000000000000000000.000000000000000000000000000000
```

and all calculations could be carried out keeping 34 digits to the left of the decimal point and 28 places to the right of it. Doing so would allow 62 significant digits in the results. On a binary computer, multiple-precision arithmetic could be used to provide enough significance. However, the mass of the sun is not even known accurately to five significant digits, let alone 62. In fact few measurements of any kind can (or need) be made accurately to 62 significant digits. Although it would be possible to keep all intermediate results to 62 significant digits and then throw away 50 or 60 of them before printing the final results, doing this is wasteful of both CPU time and memory.

What is needed is a system for representing numbers in which the range of expressible numbers is independent of the number of significant digits. In this appendix, such a system will be discussed. It is based on the scientific notation commonly used in physics, chemistry, and engineering.

## B.1 PRINCIPLES OF FLOATING POINT

One way of separating the range from the precision is to express numbers in the familiar scientific notation

$$n = f \times 10^e$$

where  $f$  is called the **fraction**, or **mantissa**, and  $e$  is a positive or negative integer called the **exponent**. The computer version of this notation is called **floating point**. Some examples of numbers expressed in this form are

$$\begin{aligned} 3.14 &= 0.314 \times 10^1 = 3.14 \times 10^0 \\ 0.000001 &= 0.1 \times 10^{-5} = 1.0 \times 10^{-6} \\ 1941 &= 0.1941 \times 10^4 = 1.941 \times 10^3 \end{aligned}$$

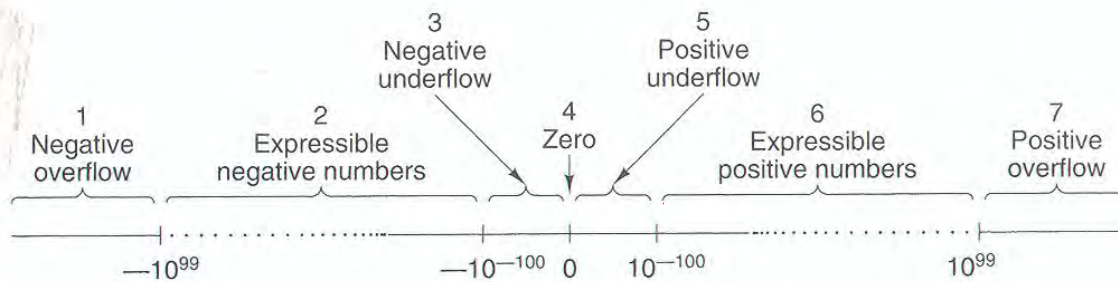
The range is effectively determined by the number of digits in the exponent and the precision is determined by the number of digits in the fraction. Because there is more than one way to represent a given number, one form is usually chosen as the standard. In order to investigate the properties of this method of representing numbers, consider a representation,  $R$ , with a signed three-digit fraction in the range  $0.1 \leq |f| < 1$  or zero and a signed two-digit exponent. These numbers range in magnitude from  $+0.100 \times 10^{-99}$  to  $+0.999 \times 10^{+99}$ , a span of nearly 199 orders of magnitude, yet only five digits and two signs are needed to store a number.

Floating-point numbers can be used to model the real-number system of mathematics, although there are some important differences. Figure B-1 gives a grossly exaggerated schematic of the real number line. The real line is divided up into seven regions:

1. Large negative numbers less than  $-0.999 \times 10^{99}$ .
2. Negative numbers between  $-0.999 \times 10^{99}$  and  $-0.100 \times 10^{-99}$ .
3. Small negative numbers with magnitudes less than  $0.100 \times 10^{-99}$ .
4. Zero.
5. Small positive numbers with magnitudes less than  $0.100 \times 10^{-99}$ .
6. Positive numbers between  $0.100 \times 10^{-99}$  and  $0.999 \times 10^{99}$ .
7. Large positive numbers greater than  $0.999 \times 10^{99}$ .

One major difference between the set of numbers representable with three fraction and two exponent digits and the real numbers is that the former cannot be used to express any numbers in regions 1, 3, 5, or 7. If the result of an arithmetic operation yields a number in regions 1 or 7—for example,  $10^{60} \times 10^{60} = 10^{120}$ —**overflow error** will occur and the answer will be incorrect. The reason is due to the finite nature of the representation for numbers and is unavoidable. Similarly,





**Figure B-1.** The real number line can be divided into seven regions.

a result in regions 3 or 5 cannot be expressed either. This situation is called **underflow error**. Underflow error is less serious than overflow error, because 0 is often a satisfactory approximation to numbers in regions 3 and 5. A bank balance of  $10^{-102}$  dollars is hardly better than a bank balance of 0.

Another important difference between floating-point numbers and real numbers is their density. Between any two real numbers,  $x$  and  $y$ , is another real number, no matter how close  $x$  is to  $y$ . This property comes from the fact that for any distinct real numbers,  $x$  and  $y$ ,  $z = (x + y)/2$  is a real number between them. The real numbers form a continuum.

Floating-point numbers, in contrast, do not form a continuum. Exactly 179,100 positive numbers can be expressed in the five-digit, two-sign system used above, 179,100 negative numbers, and 0 (which can be expressed in many ways), for a total of 358,201 numbers. Of the infinite number of real numbers between  $-10^{+100}$  and  $+0.999 \times 10^{99}$ , only 358,201 of them can be specified by this notation. They are symbolized by the dots in Fig. B-1. It is quite possible for the result of a calculation to be one of the other numbers, even though it is in region 2 or 6. For example,  $+0.100 \times 10^3$  divided by 3 cannot be expressed *exactly* in our system of representation. If the result of a calculation cannot be expressed in the number representation being used, the obvious thing to do is to use the nearest number that can be expressed. This process is called **rounding**.

The spacing between adjacent expressible numbers is not constant throughout region 2 or 6. The separation between  $+0.998 \times 10^{99}$  and  $+0.999 \times 10^{99}$  is vastly more than the separation between  $+0.998 \times 10^0$  and  $+0.999 \times 10^0$ . However, when the separation between a number and its successor is expressed as a percentage of that number, there is no systematic variation throughout region 2 or 6. In other words, the **relative error** introduced by rounding is approximately the same for small numbers as large numbers.

Although the preceding discussion was in terms of a representation system with a three-digit fraction and a two-digit exponent, the conclusions drawn are valid for other representation systems as well. Changing the number of digits in the fraction or exponent merely shifts the boundaries of regions 2 and 6 and changes the number of expressible points in them. Increasing the number of digits in the fraction increases the density of points and therefore improves the accuracy



of approximations. Increasing the number of digits in the exponent increases the size of regions 2 and 6 by shrinking regions 1, 3, 5, and 7. Figure B-2 shows the approximate boundaries of region 6 for floating-point decimal numbers for various sizes of fraction and exponent.

Digits in fraction	Digits in exponent	Lower bound	Upper bound
3	1	$10^{-12}$	$10^9$
3	2	$10^{-102}$	$10^{99}$
3	3	$10^{-1002}$	$10^{999}$
3	4	$10^{-10002}$	$10^{9999}$
4	1	$10^{-13}$	$10^9$
4	2	$10^{-103}$	$10^{99}$
4	3	$10^{-1003}$	$10^{999}$
4	4	$10^{-10003}$	$10^{9999}$
5	1	$10^{-14}$	$10^9$
5	2	$10^{-104}$	$10^{99}$
5	3	$10^{-1004}$	$10^{999}$
5	4	$10^{-10004}$	$10^{9999}$
10	3	$10^{-1009}$	$10^{999}$
20	3	$10^{-1019}$	$10^{999}$

**Figure B-2.** The approximate lower and upper bounds of expressible (unnormalized) floating-point decimal numbers.

A variation of this representation is used in computers. For efficiency, exponentiation is to base 2, 4, 8, or 16 rather than 10, in which case the fraction consists of a string of binary, base-4, octal, or hexadecimal digits. If the leftmost of these digits is zero, all the digits can be shifted one place to the left and the exponent decreased by 1, without changing the value of the number (barring underflow). A fraction with a nonzero leftmost digit is said to be **normalized**.

Normalized numbers are generally preferable to unnormalized numbers, because there is only one normalized form, whereas there are many unnormalized forms. Examples of normalized floating-point numbers are given in Fig. B-3 for two bases of exponentiation. In these examples a 16-bit fraction (including sign bit) and a 7-bit exponent using excess 64 notation are shown. The radix point is to the left of the leftmost fraction bit—that is, to the right of the exponent.