

演習問題 1. 問 2 の略解

大石進一

May 27, 2003

§1 Dekkerの定理の証明

定義 有理数 x が t ビットの浮動小数点数とは正の整数 t が存在して、

$$x = m2^k$$

と書けることをいう。ただし、 m と k は整数で、 $|m| < 2^t$ を満たしているものとする。また、正規化条件

$$2^{t-1} \leq |m|$$

を仮定する。

以上の正規化の条件の下では 2^k が最下位桁 (t bit 目) となるので、 $u = 2^k$ を $\text{ulp}(x)$ と表す。0 は正規化条件を満たすことができないので別に定義し、これも t bit 浮動小数点数とよぶことにする (例えば $0 = 0 * 2^0$ など)。 $\text{ulp}(0) = 0$ と定義しておく。マシンエプシロンは 2^{1-t} となる。

浮動小数点数演算が faithful であること

実数 r から t bit 浮動小数点数 x への丸めの演算子を

$\bigcirc r = r$ に最も近い t bit 浮動小数点数 x , 2つあるときは偶数丸めによって定義する。

このとき、 $\cdot \in \{+, -, \times, /\}$ として、 t bit 浮動小数点数 x, y に対して

$$x \odot y = \bigcirc(x \cdot y)$$

が成り立つ浮動小数点演算を faithful という。

定理 1 [Sterbenz] a と b を t bit の浮動小数点数とする。

$$\frac{1}{2} \leq \frac{a}{b} \leq 2$$

が成り立つなら、 $a - b$ は t bit 浮動小数点数となる。演算が faithful なら

$$a \ominus b = a - b$$

が成り立つ。

証明 簡単のため、 $a > b > 0$ の場合を考える。 $\text{ulp}(a) \geq \text{ulp}(b)$ であるから a と b はともに $\text{ulp}(b)$ の整数倍である。したがって、

(I) $a - b$ も $\text{ulp}(b)$ の整数倍である。

仮定から、 $a \leq 2b$ であるから

(II) $a - b \leq b$ が成り立つ。

(I), (II) から $a - b$ は t -bit 浮動小数点数であることがわかる。faithful の定義から

$$a \ominus b = \bigcirc(a - b) = a - b$$

が成り立つ。(QED)

定理 2[Dekker] a と b を t -bit 浮動小数点数で $|a| \geq |b|$ を満たすものとする。このとき、つぎのアルゴリズムによって nonoverlapping expansion $x + y$ が計算でき、 $a + b = x + y$ を満たす:

```
function [x, y] = FastTwoSum(a, b)
x = a ⊕ b;
bV = x ⊖ a;
y = b ⊖ bV;
```

証明 簡単のため、 $a > b > 0$ の場合を考える。関数 `FastTwoSum` の2行目を考える。

$$1 \leq \frac{x}{a} \leq 2$$

であるから、定理1により、

$$b_V = x \ominus a = x - a$$

となる。

ここで

$$x = a \oplus b = a + b + \text{err}(a \oplus b)$$

と書くと、

$$|\text{err}(a \oplus b)| \leq b \tag{1}$$

が示せる(少し後で証明する)。 a も b も $\text{ulp}(b)$ の整数倍であったから、 $a + b$ も $\text{ulp}(b)$ の整数倍である。 $b < a \leq a \oplus b$ であるので、 $a \oplus b$ も $\text{ulp}(b)$ の整数倍となる。よって、

$$\text{err}(a \oplus b) = (a + b) - a \oplus b$$

も $\text{ulp}(b)$ の整数倍になっている。ここで、(1)から、 $\text{err}(a \oplus b)$ は t bit浮動小数点数になることがわかる。

(1)の成立することは次の図から説明できるが、これは講義の際に詳しく行う。

