

## 5 . 多重精度算術演算<sup>1</sup>

Delphi における整数型は、8 ビット長のもの (Shortint、Byte)、16 ビット長のもの (Shortint、Word)、32 ビット長のもの (Longint) がある。Delphi 5 ではさらに、32 ビット長のものとして Longword、64 ビット長のものとして Int64 が用意されている。これらの整数型より桁数の多い整数を扱うときは、必要な桁数に対応した型を用意する必要がある。用意した型のデータの処理・演算は、多重精度算術演算<sup>(1)(2)</sup>と呼ばれているアルゴリズムによって行うことができる。多重精度算術演算を、ここでは2進数表記に対して行う。まず、有限桁数の2進数表記の演算について説明する。

### 2 進数表記の演算

正整数 9 は2進数で表わすと

$$9 = 1 \times 2^3 + 0 \times 2^2 + 0 \times 2^1 + 1 \times 2^0$$

となる。すなわち、

$$9_{10} = 1001_2$$

と表わせる。数値の右下に付けられている10や2の数は、それぞれ10進数表記、2進数表記であることを明記するためのものである。

整数 30 は

$$30 = 1 \times 2^4 + 1 \times 2^3 + 1 \times 2^2 + 1 \times 2^1 + 0 \times 2^0$$

と表わせる。すなわち、

$$30_{10} = 11110_2$$

と表わせる。

2進数4桁で0～15の数が表わせる。2進数は4桁ずつ区切ることによって簡単に16進数表記に変換できる。11110<sub>2</sub>を右から4桁ずつに区切ると、1と1110に分けられる(図1)。

---

<sup>1</sup> この解説は、TRY!PC,2000年7月号「Delphiによる数値計算プログラミングのすすめ：第5回 多重精度算術演算・方程式の根の計算」の原稿をもとにしたものです。

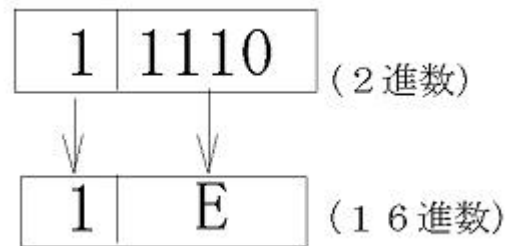


図1 2進数の16進数への変換。Eは16進数において14を表わす。

1110は10進数で $8+4+2=14$ を表わすが、14は16進数では英字eまたはEで表わす。したがって、

$$30_{10} = 11110_2 = 1E_{16}$$

と書ける。16進数は、Delphiでは\$を付けて\$1Eと表わし、C++では先頭に0xまたは0Xを付けて0X1Eと表わす。

コンピュータのメモリーは、8ビットを1バイトとして数える。1ビットは2進数の1桁に対応する。4ビットで2進数4桁、すなわち16進数1桁、を表わすので、1バイトは16進数で2桁を表わすことになる。

コンピュータにおいて、整数は、必要ならば左側に0を置くことにより、型によって決まった桁数（ビット数）の2進数で表わされる。Shortint型の場合は1バイトすなわち8ビットで表わされるので、2進数で8桁までの数が表わされる。演算の結果が8桁を超えたときは、その超えた桁はメモリーに保持されずに消える<sup>(3)</sup>。

例えば、

$$(01000000_2) \times (00001000_2) = 1000000000_2$$

の場合、右辺は8桁を超える数なので0になる（図2）。

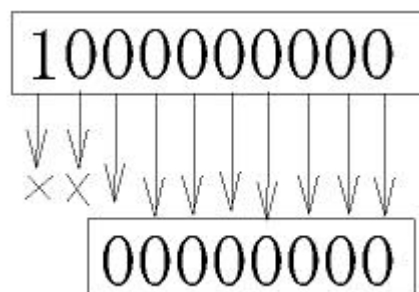


図2 Shortintは8ビットなので、演算結果が2進数で8桁を超えた部分は消える。

2進数で8桁を超えた部分が消えてしまうということは、 $2^8$ を法とする演算を行うことと同じである。すなわち、

$$(0100000)_2 \times (00001000)_2 \equiv 0 \pmod{2^8}$$

ということである。

$2^8$ を法とする演算では、

$$(00000001)_2 + (11111111)_2 \equiv 0 \pmod{2^8}$$

となる。

$2^8$ を法とする演算、すなわち Shortint 型での演算では、1に $(11111111)_2$ を足すと0になる。したがって、 $(11111111)_2$ は-1と同じ働きをしている。Shortint 型では、最高位の桁が1である数は、それとの加算を行ったとき0となる数の負の値を表わすとみなす。2の負の数 -2は、

$$(00000010)_2 + (11111110)_2 = 100000000_2 \equiv 0 \pmod{2^8}$$

となるので、 $(11111110)_2$ が -2を表わす。

コンピュータでは、数値は2進数で表わされ、演算は2のべき乗を法とする演算になっている。 $n$ ビット( $n$ 桁)の表記では、 $2^n$ を法とする演算になる。負の数は、最高位の桁が1のもので表わされる。

$n = 8$ の場合の $2^n$ を法とする演算結果を確認するプロジェクトを PCkModular.dpr として用意した。このプロジェクトを実行すると、図3のようなフォームが表示される。

図 3 PckModular.dpr の実行開始時のフォーム

上の行と中の行に 2 つの整数値を設定する。設定は、Byte と ShortInt の 2 つの欄のいずれかで行う。Binary の表示の下欄の Edit コンポーネントは表示専用である。Byte の下の欄は、8 桁までの 2 進数を 0 ~ 255 の範囲で設定する。この範囲を超えたり、不適当な文字が設定されると、範囲内の数値に適当に変換されたり、空白になったりする。Byte の欄に値を設定すると、そのときの値に対応した 2 進数表記が Binary の欄に表示される。ShortInt の欄には、8 桁目を符号ビットとみた符号付き整数としての値が表示される。例えば、Byte 欄に 255 を設定すると、ShortInt 欄には - 1 が表示される。逆に、ShortInt 欄に - 1 を設定すると、Byte 欄には 255 が表示される。

上段と中段にそれぞれ数値を設定してから、「Sum」あるいは「Mul」ボタンをクリックすると、上段の数と中段の数の和、あるいは、積が算出されて下段に表示される。 - 1 と - 128 の和が図 4 に、積が図 5 に示されている。

Byte	Binary	ShortInt
255	11111111	-1
128	10000000	-128
Result	Result(Binary)	Result(ShortInt)
127	01111111	127

Buttons: Sum, Mul, Exit

Callout: Sum」のクリックで和の計算

図4  $2^8$ を法とする和の計算

Byte	Binary	ShortInt
255	11111111	-1
128	10000000	-128
Result	Result(Binary)	Result(ShortInt)
128	10000000	127

Buttons: Sum, Mul, Exit

Callout: Mul」のクリックで積の計算

図5  $2^8$ を法とする積の計算

普通の10進数での演算結果が8桁を超えるときは、このように普通の整数値の計算とは異なった感じの結果になることがある。しかし、普通の10進数での演算結果が8桁を超えない場合、例えば - 1 と - 2 の計算の場合は、図6 a、図6 bのように ShortInt の欄の値は普通の整数値の演算と同じ結果になっている。図6の場合、Byte 欄および Binary 欄に

おける演算結果は  $2^8$  を法とするものである。

Byte	Binary	ShortInt
255	11111111	-1
254	11111110	-2
Result	Result(Binary)	Result(ShortInt)
253	11111101	-3

Sum Mul Exit

図 6 a - 1 と - 2 の和

Byte	Binary	ShortInt
255	11111111	-1
254	11111110	-2
Result	Result(Binary)	Result(ShortInt)
2	00000010	2

Sum Mul Exit

図 6 b - 1 と - 2 の積

### Knuth(1998)<sup>(2)</sup> の表記法

正整数  $p$  の 2 進数表記

$$p = u_{n-1} \times 2^{n-1} + u_{n-2} \times 2^{n-2} + \cdots + u_1 \times 2^1 + u_0 \times 2^0$$

をここでは Knuth ( 1998 ) に従って

$$p = (u_{n-1}u_{n-2} \cdots u_1u_0)_2$$

と表わす。Knuth(1981)の訳 (1986)<sup>(1)</sup> では、添字は  $(u_1u_2 \cdots u_n)_2$  というように昇順であるが、ここでは Knuth(1998)の方に従う。

正整数  $p$  が  $b$  進法で

$$p = v_{m-1} \times b^{m-1} + v_{m-2} \times b^{m-2} + \cdots + v_1 \times b^1 + v_0 \times b^0$$

と表わされるときは

$$p = (v_{m-1}v_{m-2} \cdots v_1v_0)_b$$

と表わす。  $b$  は基数 (base または radix) という。

コンピュータでは、2 進数 4 桁をまとめて、 $2^4 = 16$  を基数として整数を表わすことがある。ここでは、プログラミングが簡単になるように、2 進数 7 桁をまとめて、 $2^7 = 128$  を基数とする。128 を基数とする 15 桁の表記法は、2 進数で  $7 \times 15 = 105$  桁の数が扱える。 $2^{105}$  は 10 進数で 32 桁の数なので、128 を基数とする 15 桁の表記法では、10 進数で 32 桁ぐらいまでの整数が扱えることになる。

$n$  ビットの整数が 10 進数による表記で何桁になるのかを求めるプロジェクトとして PCheck.dpr を用意した。プロジェクト PCheck.dpr を実行すると、図 7 のようなフォームが表示される。

図 7 ビット数から 10 進数での桁数を求める

「ビット数」の右側の Edit コンポーネント内にビット数を設定して「Calc」ボタンをクリックすると、「桁数」の右側に 10 進数表記での桁数が表示される。

### 多重精度整数型 TBigInt

多重精度算術演算のための型 TBigInt をユニットファイル UBigInt.pas に用意した。TBigInt は、以下のように宣言されている。

```
type TBigInt = array[0..NBigInt-1] of byte;
```

TBigInt 型は、byte 型を要素とする配列である。byte 型の各要素の最上位のビット (MSB、Most Significant Bit) はキャリービットとして使い、1 byte 中の下位 7 ビットで 2 進数の各桁を表わす。したがって、byte 型の各要素は、 $2^7 = 128$  進数表記での各桁に対応する (図 8)。

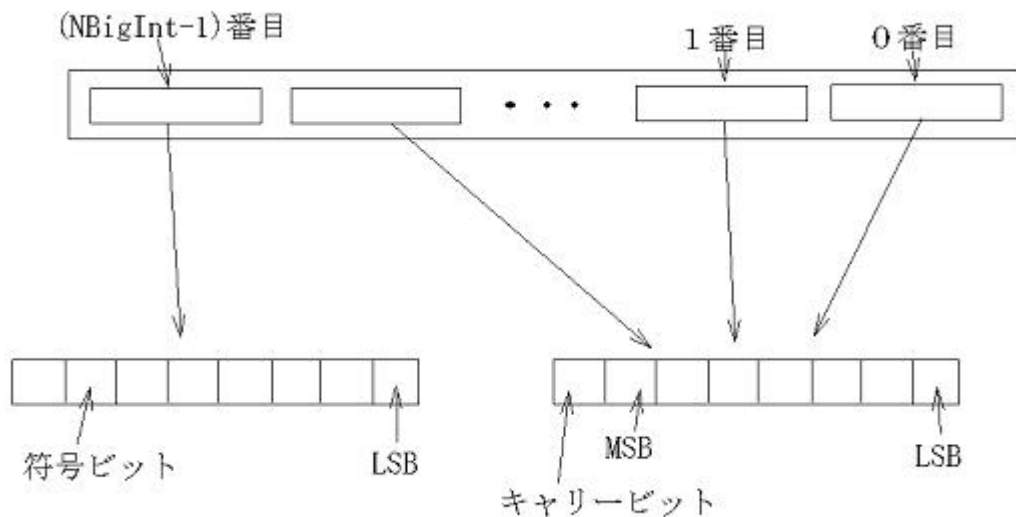


図 8 TBigInt 型の構成

ユニットファイル UBigInt.pas では

```
const NBigInt = 15;
```

と宣言されているので、 $7 \times 15 = 105$  ビットの整数が扱える。105 ビットの中で最上位ビットは符号ビットとして使うので、 $-2^{104}$  から  $2^{104} - 1$  までの整数値を表わすことができる。

プロジェクト PTable.dpr は、TBigInt 型を用いて多重精度の整数値の和を求めるものである。整数値はグリッドのセルに設定する (図 9)。





図9 和を求める整数値の設定

グリッドは、「追加」ボタンのクリックで行を追加することができる。データ設定用セルをクリックしてから「追加」ボタンをクリックすると、クリックしておいたセル（アクティブなセル）の下に行が追加される。不要なセルは、そのセルをクリックしてから「削除」ボタンをクリックすると削除される。全てのデータ設定用セルに数値を設定した後、「計算」ボタンをクリックすると、データの総和が下の Edit コンポーネントに表示される（図 10）。

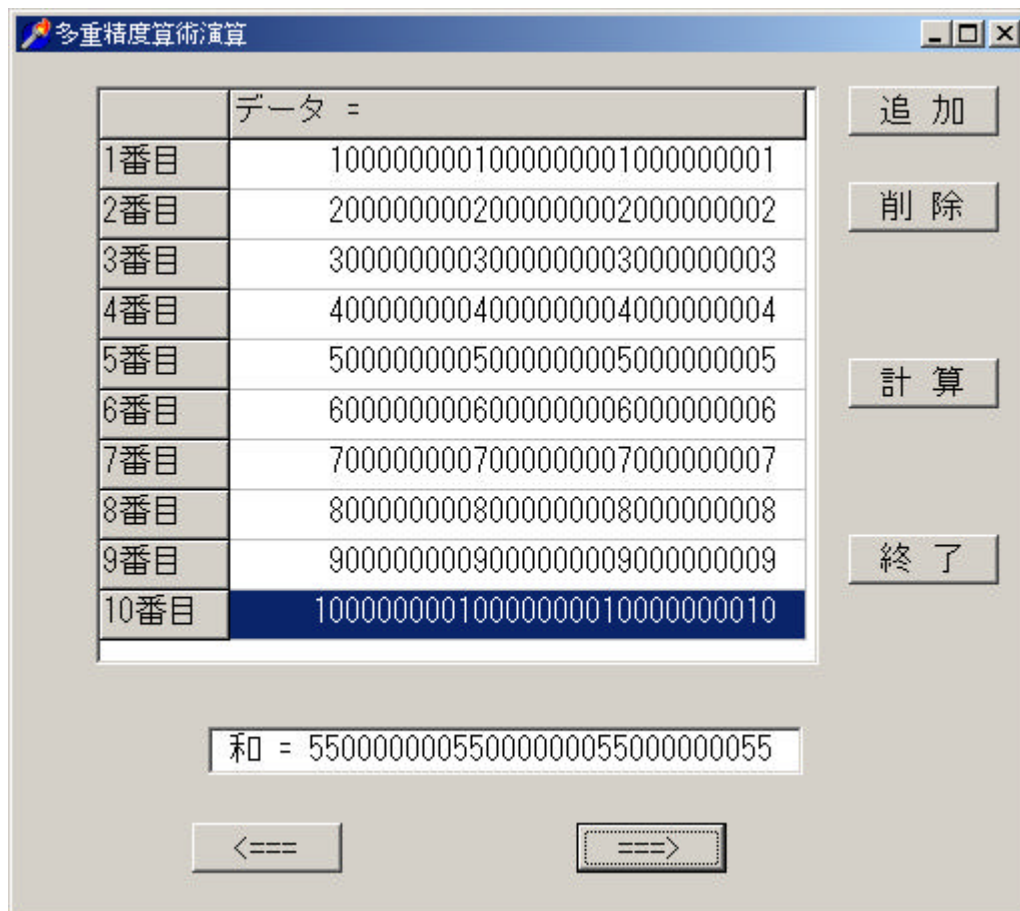


図 10 「計算」ボタンのクリックで和が求まる

「===>」ボタンをクリックすると、表示が右詰になる。左詰にするときは、「<===」ボタンをクリックする。データの訂正のときは、左詰の方がやり易い。

### ユニット UBigInt

多重精度演算のためのユニットファイル UBigInt.pas には、表 1 の関数が用意されている。

表1 ユニット UBigInt.pas に用意されている関数

function BigAdd( a, b : TBigInt ) : TBigInt;
a と b の和を関数値とする
function BigSub( a, b : TBigInt ) : TBigInt;
a と b の差を関数値とする
function BigMul( a, b : TBigInt ) : TBigInt;
a と b の積を関数値とする
function BigDiv( a, b : TBigInt ) : TBigInt;
a を b で割った商を関数値とする
function BigMod( a, b : TBigInt ) : TBigInt;
a を b で割った剰余を関数値とする
function StrToBig( s : string ) : TBigInt;
文字列 s の表わす整数値を TBigInt 型として返す
function BigToStr( a : TBigInt ) : string;
a の値を表わす整数値を ( 符号と ) 数字の文字列として返す。

function BigToFloat( a : TBigInt ) : Extended;
a の表わす整数値を Extended 型として返す
function BigALB( a , b : TBigInt ) : Boolean;
a < b のときは True、それ以外のとき False を関数値とする
function BigALEB( a, b : TBigInt ) : Boolean;
a ≤ b のときは True、それ以外のときは False を関数値とする
function BigAGB( a , b : TBigInt ) : Boolean;
a > b のときは True、それ以外のときは False を関数値とする
function BigAGEB( a, b : TBigInt ) : Boolean;
a ≥ b のときは True、それ以外のときは False を関数値とする
function BigAEB( a, b : TBigInt ) : Boolean;
a = b のときは True、それ以外のときは False を関数値とする

これらの関数では、TBigInt 型の整数の加減乗除の演算を基礎として処理が行われている。これらの加減乗除の演算は、Knuth ( 1998 )<sup>( 2 )</sup> のアルゴリズムに従って行われているが、基本的には紙の上で行う筆算と同じものである。以下の説明では、TBigInt 型に合わせて、基数 b ( base、radix ) を  $2^7 = 128$  とする b 進数表記で整数 p を表わす。

$$p = u_{n-1} \times b^{n-1} + \cdots + u_0 \times b^0 = (u_{n-1} \cdots u_0)_b = (u_{n-1} \cdots u_0)_{128}$$

型 TBigInt は

```
type TBigInt = array[0..NBigInt-1] of byte;
```

と宣言されていて、各 byte 型の要素で各  $u_i$  を表わす ( 図 8 )。

```
var a : TBigInt;
```

のとき、 $a[NBigInt-1]$  が最上位桁  $u_{n-1}$  ,  $a[0]$  が最下位桁  $u_0$  を表わす。ただし、

```
n = NBigInt
```

である。

2 数、 $u = (u_{n-1} \cdots u_0)_{128}$  と  $v = (v_{n-1} \cdots v_0)_{128}$ 、の四則演算について説明する。

### 加算のアルゴリズム

$u = (u_{n-1} \cdots u_0)_{128}$  と  $v = (v_{n-1} \cdots v_0)_{128}$  の加算を次のように行う ( Knuth,1998, Algorithm A )。ここで、和を  $w = (w_{n-1} \cdots w_0)_{128}$  に求めるものとする ( 図 1 1 )。

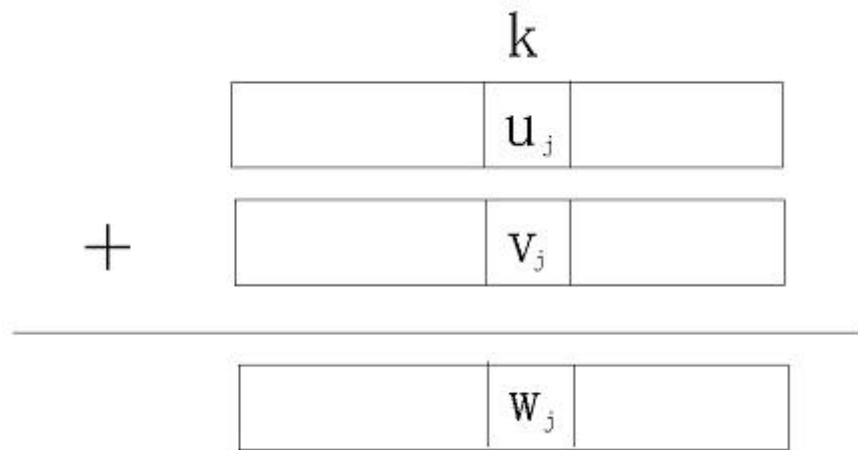


図 1 1  $u = (u_{n-1}, \cdots, u_0)_{128}$  と  $v = (v_{n-1}, \cdots, v_0)_{128}$  の足し算。  
 $k$  は桁上がりの値を表わす。

( A1 )  $j \leftarrow 0$ 、 $k \leftarrow 0$  とおく。ここで、矢印  $\leftarrow$  は代入を表わす。  $j$  は現在計算を行っている桁を表わし、 $k$  はキャリー( 桁上がり )を表わす。これらは初期値として 0 とおく。

( A2 )  $w_j = (u_j + v_j + k) \bmod 128$ 、

$k \leftarrow (u_j + v_j + k) \div 128$

とおく。

(A3)  $j \leftarrow j+1$  とおく。

$j < n$  ならば、(A2) に戻る。

$j \geq n$  ならば、このアルゴリズムを終了する。アルゴリズムの終了後、

$w = (w_{n-1} \cdots w_0)_{128}$  には  $u$  と  $v$  の和が格納されている。

関数 BigAdd は、上のアルゴリズムによって和を求めるものである。

### 減算のアルゴリズム

関数 BigSub では、 $a$  と  $b$  の差を  $a$  と  $-b$  の和として求めている。

負の数は、 $2^{105}$  を法とする演算においては、最上位のビットが 1 のもので表わされる。したがって、 $-b$  は  $b$  を表わすビットを反転 (1 は 0、0 は 1 にする) してから 1 を加えることによって得られる。

$a$  と  $-b$  の和は、関数 BigAdd によって求めている。

### 乗算のアルゴリズム

2 数、 $u = (u_{n-1} \cdots u_0)_{128}$  と  $v = (v_{n-1} \cdots v_0)_{128}$ 、の積  $w = (w_{2n-1} \cdots w_0)_{128}$  は、以下の手順で求めている (Knuth, 1998, Algorithm M)。

(M1) まず、 $j = n-1, \dots, 0$  に対して  $w_j \leftarrow 0$ ; とおく。

次に、 $j \leftarrow 0$  とおく。 $j$  は、 $v$  の  $128^j$  の位  $v_j$  による掛け算であることを表わすものである。

(M2)  $v_j = 0$  ならば、 $w_{j+n} \leftarrow 0$  とおき、(M6) に跳ぶ。

(M3)  $i \leftarrow 0$ 、 $k \leftarrow 0$  とおく。ここで、 $i$  は  $u$  の  $128^i$  の位  $u_i$  との掛け算であることを表わし、 $k$  は桁上がりの数を表わすものである。これらは初期値として 0 とおく。

(M4)  $t \leftarrow u_i \times v_j + w_{i+j} + k$  とおく。ここで、 $t$  は例えば Word 型の変数として、桁溢れがないようにする。

$$w_{i+j} \leftarrow t \bmod 128$$

$$k \leftarrow t \operatorname{div} 128$$

とおく (図 1 2 )。

(M5)  $i \leftarrow i+1$  とおく。

$i < n$  ならば、(M4) に戻る。

$i \geq n$  ならば、 $w_{j+n} \leftarrow k$  とおく。

(M6)  $j \leftarrow j+1$  とおく。

$j < n$  ならば、(M2) に戻る。

$j \geq n$  ならば、このアルゴリズムを終了する。アルゴリズムの終了後、

$w = (w_{2n-1} \cdots w_0)_{128}$  には、 $u$  と  $v$  の積が格納されている。

関数 BigMul は、上のアルゴリズムによって積を求めるものである。

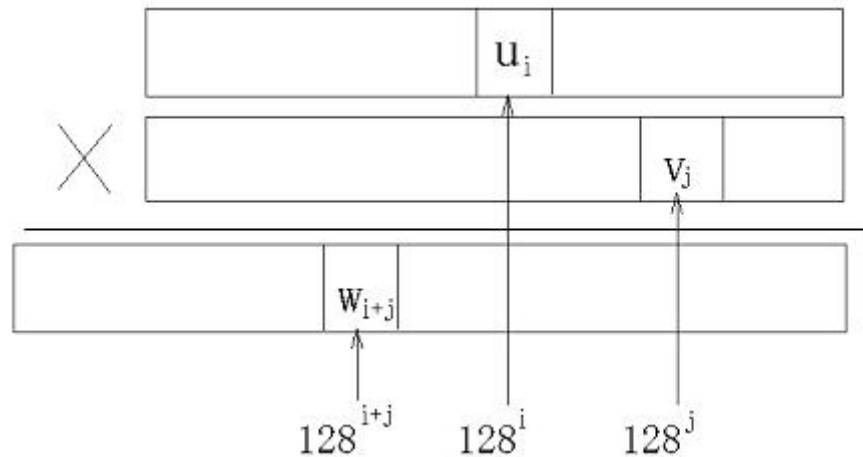


図 1 2  $u = (u_{n-1} \cdots u_0)_{128}$  と  $v = (v_{n-1} \cdots v_0)_{128}$  の掛け算において、 $u$  の  $128^i$  の位  $u_i$  と  $v$  の  $128^j$  の位  $v_j$  の積は  $w = u \times v$  の  $128^{i+j}$  の位  $w^{i+j}$  に加えられる

### 除算のアルゴリズム

2 数、 $u = (u_{n-1} \cdots u_0)_{128}$  と  $v = (v_{n-1} \cdots v_0)_{128}$ 、の商  $q = (q_{n-1} \cdots q_0)_{128}$  および剰余  $r = (r_{n-1} \cdots r_0)_{128}$  の算出 (Knuth、1998、Algorithm D) においては、まず、 $u$  および  $v$  が非負であるようにしてから商と剰余を求める。 $u$  あるいは  $v$  が負のときは、非負のときの商と剰余から負のときの商と剰余を求める。例えば、 $u < 0$ 、 $v > 0$  のときは、

$$-u = vq + r$$

から

$$u = v(-q) + (-r)$$

として商  $q$  と剰余  $r$  を求める。

以下の説明では、 $u$  および  $v$  は非負であるとする。

$u < v$  のときは、 $u = v \times 0 + u$  となり、 $q = 0$ ,  $r = u$  となる。

$u \geq v$  の場合について考える。

$v_j = 0$ ,  $j = n-1, \dots, 1$  のときは  $v_0$  で  $u_i$  を順番に割っていけばよい。以下では、ある

$j > 0$  に対して  $v_j \neq 0$  である場合について考える。

$v_j > 0$  である最大の  $j$  を  $n-1$  とおく。このとき、適当な整数  $m$  に対して、

$u = (u_{m+n} \cdots u_0)_{128}$ 、 $v = (v_{n-1} \cdots v_0)_{128}$  と表わすことができる。

筆算による割り算は、上の桁から計算していく (図 13)。

$$\begin{array}{r} q_m \cdots q_0 \\ v_{n-1} \cdots v_0 \overline{) u_{m+n} \cdots u_m \cdots u_0} \\ \cdot \cdot \cdot \end{array}$$

図 13  $(u_{m+n} \cdots u_0)_{128} \div (v_{n-1} \cdots v_0)$  の筆算

最初の桁  $q_m$  の計算について考える。

$$(u_{m+n} \cdots u_m)_{128} \div (v_{n-1} \cdots v_0)_{128} < 128 \quad (1)$$

であるとする。この商が 128 以上のときは、被除数  $u$  の先頭に 0 を付け加えておく。すなわち、 $u_{m+n+1} = 0$ 、 $m \leftarrow m+1$  として、 $(0u_{m+n} \cdots u_0)_{128}$  を改めて  $(u_{m+n}u_{m+n-1} \cdots u_0)_{128}$  とおくことによって商が 128 より小さくなるようにする。いま、 $\hat{q}$  を次式によって定める。

$$\hat{q} = \min\{(u_{m+n} \times 128 + u_{m+n-1}) \operatorname{div} v_{n-1}, 128-1\} \quad (2)$$

このとき、

$$q_m = (u_{m+n} \cdots u_m)_{128} \operatorname{div} (v_{n-1} \cdots v_0)_{128}$$

とおくと

$$\hat{q} \geq q_m \quad (3)$$

が成り立つ (Knuth, 1998, p.271, Theorem A)。式 (3) の証明は、附 A に示した。

さらに、



$$v_{n-1} \geq 128 \div 2 = 64 \quad (4)$$

であれば、

$$q_m \leq \hat{q} \leq q_m + 2 \quad (5)$$

となる (Knuth, 1998, p.271, Theorem B)。式 (5) の証明は附 B に示した。

(5) 式より、商  $q_m$  は  $\hat{q}$ 、 $\hat{q}-1$ 、 $\hat{q}-2$  の 3 つの値の 1 つに一致することがわかる。 $q_m$  を求めるときは、 $\hat{q}$  を商  $q_m$  の最初の推定値とする。この推定値が商として大き過ぎるときは適切な値になるまで 1 つずつ減らしていく。この減らす手続きは、式 (5) より高々 2 回で済むことになる。商  $q_m$  の見当をつけるのに、被除数の上 2 桁と除数の上 1 桁の商  $\hat{q}$  で行うのが式 (2) である。5 4 3 2 1  $\div$  9 8 の最初の桁を 5 4  $\div$  9 によって見当をつけるということである。このとき、除数の上 1 桁の値が小さいと 2 桁目の値の影響を受け易いので、基数の 1/2 より大きくしておくというのが (4) 式の条件である。

$$v_{n-1} < 64$$

のときは、次のように調整することによって (4) 式が成り立つようにすることができる。

一般に、 $u$  を  $v$  で割ったときの商を  $q$ 、剰余を  $r$  とおくと、

$$u = vq + r$$

が成り立つ。両辺に  $a$  をかけると、

$$(au) = (av)q + (ar)$$

が成り立つ。すなわち、除数と被除数を定数倍しても、商は変わらない。このことより、 $v$  が 64 より小さいときは、 $u$  と  $v$  を定数倍することにより (4) 式 ( $v \geq 64$ ) が成り立つようにして、商を求めることができる。

上のことを踏まえて、 $n > 1$  の場合の割り算を以下のように行う (Knuth, 1998, Algorithm D)。

$u = (u_{m+n} \cdots u_0)_{128}$  を  $v = (v_{n-1} \cdots v_0)_{128}$  で割ったときの商を  $q = (q_m \cdots q_0)_{128}$ 、剰余を  $r = (r_{n-1} \cdots r_0)_{128}$  とおく。

(D 1) 上の式 (5) の説明での条件「 $v_{n-1} \geq 64$ 」が成り立つように、

$u = (u_{m+n} \cdots u_0)_{128}$  と  $v = (v_{n-1} \cdots v_0)_{128}$  の形を整えておく。 $u$  と  $v$  の

定数倍に用いた定数を  $d$  とおく。

( D 2 )  $j \leftarrow m$  とおく。ここで、 $j$  は商の  $128^j$  の位  $q_j$  の計算であること  
を表わすものである。

( D 3 )  $\hat{q} \leftarrow \{(u_{j+n} \times 128 + u_{j+n-1}) \operatorname{div} v_{n-1}\}$

$\hat{r} \leftarrow \{(u_{j+n} \times 128 + u_{j+n-1}) \bmod v_{n-1}\}$

とおく。

( D 3 a )  $\hat{q} \geq 128$  または  $\hat{q}v_{n-2} > 128 \times \hat{r} + u_{j+n-2}$

ならば、

$\hat{q} \leftarrow \hat{q} - 1, \hat{r} \leftarrow \hat{r} + v_{n-1}$

とおく。

条件「 $\hat{q} \geq 128$  または  $\hat{q}v_{n-2} > 128 \times \hat{r} + u_{j+n-2}$ 」が成り

立たないときは、( D 4 ) に進む。

( D 3 b ) ( D 3 a ) に戻る。

( D 4 )  $(u_{j+n} \cdots u_j)_{128} \leftarrow \{(u_{j+n} \cdots u_j)_{128} - \hat{q} \times (v_{n-1} \cdots v_0)_{128}\}$  とおく。

( D 5 )  $q_j \leftarrow \hat{q}$  とおく。

$(u_{j+n} \cdots u_j)_{128} < 0$  のときは、( D 6 ) に進む。

$(u_{j+n} \cdots u_j)_{128} \geq 0$  のときは、( D 7 ) に跳ぶ。

( D 6 )  $q_j \leftarrow q_j - 1$

$(u_{j+n} \cdots u_j)_{128} \leftarrow (u_{j+n} \cdots u_j)_{128} + (v_{n-1} \cdots v_0)_{128}$

とおく。

( D 7 )  $j \leftarrow j - 1$  とおく。

$j \geq 0$  ならば、( D 3 ) に戻る。

$j < 0$  ならば、( D 8 ) に進む。

( D 8 )  $(q_m \cdots q_0)_{128}$  に商が求まっている。

剰余は、( D 1 ) での定数  $d$  で割ったもの、

$$(r_{n-1} \cdots r_0) \leftarrow \{(u_{n-1} \cdots u_0) \operatorname{div} d\}$$

として得ることができる。

上のアルゴリズムによって商と剰余を求める手続きとして BigDivMod を用意した。商を求める関数 BigDiv、および剰余を求める関数 BigMod では、手続き BigDivMod を呼出している。

なお、上のアルゴリズムの (D3a) は、以下の性質に基づいている。

いま、

$$(u_{j+n} \cdots u_j)_{128} = q \times (v_{n-1} \cdots v_0)_{128} + r, \quad 0 \leq r < (v_{n-1} \cdots v_0)_{128}$$

$$u_{j+n} \times 128 + u_{j+n-1} = \hat{q} v_{n-1} + \hat{r}$$

とおく。ここで、 $\hat{r} < 128$  の制約はない。

このとき、次の (A) および (B) が成り立つ。

(A)  $\hat{q} v_{n-2} > 128 \hat{r} + u_{j+n-2}$  ならば、 $q < \hat{q}$  (Knuth, 1998, p.282, Exercise 19)

(B)  $\hat{q} v_{n-2} \leq 128 \hat{r} + u_{j+n-2}$  かつ  $\hat{q} < 128$  ならば、 $\hat{q} \leq q + 1$  (Knuth, 1998, p.282, Exercise 20)

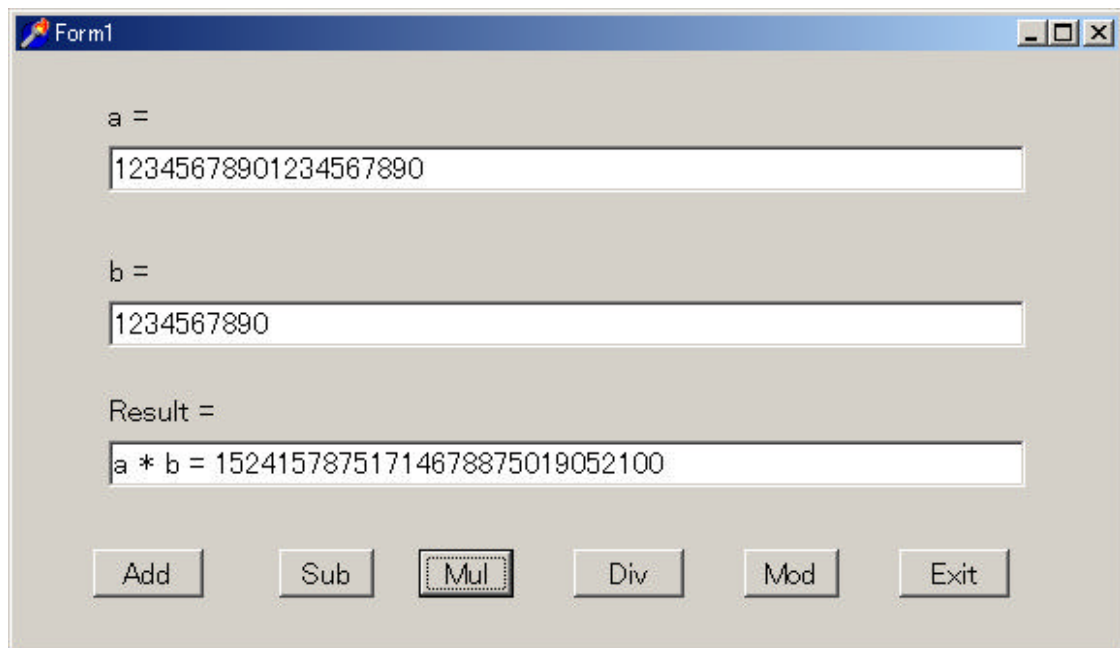
性質 (A) は、 $\hat{q}$  が  $q$  より大きくなる場合の条件を示している。なお、(3) 式より  $\hat{q}$  は  $q$  以上の数である。

性質 (B) は、 $\hat{q}$  と  $q$  の差が高々 1 である場合の条件を示している。この性質 (B) はステップ (D3) で用いている。

なお、Knuth(1998)の Algorithm D (p.272) におけるステップ D3 では、条件  $\hat{r} < 128$  のチェックが行われている。しかし、性質 (B) が成り立つので、手続き BigDivMod では  $\hat{r} < 128$  のチェックを省いている。

#### 四則演算のチェック

プロジェクト PCalc.dpr は、上の四則演算の関数を用いて加減乗除の結果を表示するものである。実行時のフォームは図 14 のようになっている。



Form1

a =  
12345678901234567890

b =  
1234567890

Result =  
a \* b = 15241578751714678875019052100

Add Sub **Mul** Div Mod Exit

図 1 4 四則演算の実行例

それぞれ「a = 」と「b = 」の下に Edit コンポーネントのところに整数値を設定して、「Add」、「Sub」、「Mul」、「Div」、「Mod」のボタンをクリックすると、a と b の和、差、積、商、剰余が「Result」の下に Edit コンポーネントに表示される。

## 附 A . 式 ( 3 ) の証明

条件 ( 1 ) のもとでは

$$q_m \leq 128 - 1 \quad (\text{A.1})$$

が成り立つ。

また、式 ( 2 ) より

$$\hat{q} \leq 128 - 1$$

が成り立つ。

$\hat{q}$  の値の場合に分けて考える。

$\hat{q} = 128 - 1$  のときは、式 ( A. 1 ) より次式

$$\hat{q} \geq q_m$$

すなわち ( 3 ) 式が成り立つ。

その他の場合、すなわち

$$\hat{q} < 128 - 1 \quad (\text{A.2})$$

の場合について考える。

式 ( 2 ) と ( A.2 ) より

$$\hat{q} = \left\lfloor \frac{u_{m+n} \times 128 + u_{m+n-1}}{v_{n-1}} \right\rfloor$$

であることがわかる。ここで記号  $[z]$  は、 $z$  を超えない最大の整数値を表わす。

したがって、

$$\hat{q} \leq \frac{u_{m+n} \times 128 + u_{m+n-1}}{v_{n-1}} \leq \hat{q} + 1 - \frac{1}{v_{n-1}}$$

が成り立つ。上式の右側の不等式より

$$\begin{aligned} u_{m+n} \times 128 + u_{m+n-1} &\leq \hat{q} v_{n-1} + v_{n-1} - 1 \\ \hat{q} v_{n-1} &\geq u_{m+n} \times 128 + u_{m+n-1} - v_{n-1} + 1 \end{aligned} \quad (\text{A.3})$$

となる。いま、

$$u^{(m)} = (u_{m+n} \cdots u_m)_{128}$$

とくと、

$$v = (v_{n-1} \cdots v_0)_{128}$$

なので、

$$\begin{aligned} u^{(m)} - \hat{q}v &= u^{(m)} - \hat{q} \cdot (v_{n-1} \cdots v_0)_{128} \\ &= u^{(m)} - \hat{q}(v_{n-1} \times 128^{n-1} + \cdots + v_0) \\ &\leq u^{(m)} - \hat{q}v_{n-1} \times 128^{n-1} \end{aligned} \quad (\text{A.4})$$

となる。

式 (A.3) と (A.4) および  $u^{(m)} = (u_{m+n} \cdots u_m)_{128}$  に注意して、

$$\begin{aligned} u^{(m)} - \hat{q}v &\leq u^{(m)} - (u_{m+n} \times 128 + u_{m+n-1} - v_{n-1} + 1) \times 128^{n-1} \\ &= u_{m+n} \times 128^n + u_{m+n-1} \times 128^{n-1} + u_{m+n-2} \times 128^{n-2} + \cdots + u_m \\ &\quad - (u_{m+n} \times 128^n + u_{m+n-1} \times 128^{n-1} - v_{n-1} \times 128^{n-1} + 128^{n-1}) \\ &= u_{m+n-2} \times 128^{n-2} + \cdots + u_m + v_{n-1} \times 128^{n-1} - 128^{n-1} \\ &= u_{m+n-2} \times 128^{n-2} + \cdots + u_m - 128^{n-1} + v_{n-1} \times 128^{n-1} \end{aligned} \quad (\text{A.5})$$

を得る。

$128^{n-1}$  は 128 進数の  $n$  桁目の位なので、128 進数表記において  $n-1$  桁までで表わされる数より大きい、すなわち

$$128^{n-1} > u_{m+n-2} \times 128^{n-2} + \cdots + u_m \quad (\text{A.6})$$

が成り立つ。

式 (A.5) と (A.6) より

$$\begin{aligned} u^{(m)} - \hat{q}v &< v_{n-1} \times 128^{n-1} \\ &\leq v_{n-1} \times 128^{n-1} + v_{n-2} \times 128^{n-2} + \cdots + v_0 \\ &= v \end{aligned}$$

となる。すなわち、

$$u^{(m)} < \hat{q}v + v = (\hat{q} + 1)v$$

となる。上式は、 $u^{(m)}$  を  $v$  で割ったときの商  $q_m$  が  $\hat{q}+1$  より小さいこと

$$q_m < \hat{q} + 1$$

を表わしている。 $q_m$  と  $\hat{q}$  は整数なので、上の不等式より

$$q_m \leq \hat{q}$$

が導ける。上式は ( 3 ) 式と同じ内容である。

## 附 B . 式 ( 5 ) の証明

( 3 ) 式が成り立つので、( 5 ) 式が成り立つことを示すためには、( 4 ) 式の仮定

$$v_{n-1} \geq 64$$

のもとで、次式

$$\hat{q} \leq q_m + 2 \quad (\text{B.1})$$

が成り立つことを示せばよい。

背理法で証明する。

次式

$$\hat{q} \geq q_m + 3 \quad (\text{B.2})$$

が成り立つとする。

( 2 ) 式より

$$\begin{aligned} \hat{q} &\leq \frac{u_{m+n} \times 128 + u_{m+n-1}}{v_{n-1}} \\ &= \frac{u_{m+n} \times 128^n + u_{m+n-1} \times 128^{n-1}}{v_{n-1} \times 128^{n-1}} \\ &\leq \frac{u_{m+n} \times 128^n + u_{m+n-1} \times 128^{n-1} + u_{m+n-2} \times 128^{n-2} + \cdots + u_m}{v_{n-1} \times 128^{n-1}} \\ &= \frac{u^{(m)}}{v_{n-1} \times 128^{n-1}} \end{aligned} \quad (\text{B.3})$$

が導かれる。 $u^{(m)}$  は附. A で定義された  $u^{(m)} = (u_{m+n} \cdots u_m)_{128}$  のことである。

また、

$$v = v_{n-1} \times 128^{n-1} + v_{n-2} \times 128^{n-2} + \cdots + v_0$$

$$< v_{n-1} \times 128^{n-1} + 128^{n-1}$$

すなわち、

$$v_{n-1} \times 128^{n-1} > v - 128^{n-1} \quad (\text{B.4})$$



が成り立つ。

式 (B.3) と (B.4) より

$$\hat{q} < \frac{u^{(m)}}{v - 128^{n-1}} \quad (\text{B.5})$$

となる。

また、

$$q_m = \left\lceil \frac{u^{(m)}}{v} \right\rceil$$

なので

$$q_m > \frac{u^{(m)}}{v} - 1 \quad (\text{B.6})$$

となる。

式 (B.2) と (B.6) より

$$3 \leq \hat{q} - q_m < \hat{q} - \left( \frac{u^{(m)}}{v} - 1 \right)$$

となる。上式と式 (B.5) より

$$\begin{aligned} 3 &< \frac{u^{(m)}}{v - 128^{n-1}} - \left( \frac{u^{(m)}}{v} - 1 \right) \\ &= \frac{u^{(m)}v - u^{(m)}(v - 128^{n-1})}{(v - 128^{n-1})v} + 1 \\ &= \frac{u^{(m)}}{v} \cdot \frac{128^{n-1}}{v - 128^{n-1}} + 1 \end{aligned}$$

すなわち、

$$\begin{aligned} 3 - 1 &< \frac{u^{(m)}}{v} \cdot \frac{128^{n-1}}{v - 128^{n-1}} \\ \frac{u^{(m)}}{v} &> 2 \cdot \frac{v - 128^{n-1}}{128^{n-1}} \\ &= 2 \cdot \frac{v_{n-1} \times 128^{n-1} + v_{n-2} \times 128^{n-2} + \dots + v_0 - 128^{n-1}}{128^{n-1}} \\ &\geq 2 \cdot \frac{v_{n-1} \times 128^{n-1} - 128^{n-1}}{128^{n-1}} \\ &= 2(v_{n-1} - 1) \quad (\text{B.7}) \end{aligned}$$

となる。

式 ( 2 ) より

$$128-1 \geq \hat{q}$$

が成り立つので、

$$\begin{aligned} 128-4 &\geq \hat{q}-3 \\ &\geq q_m \end{aligned} \quad (\text{式 (B.2) より})$$

$$\begin{aligned} &= \left\lceil \frac{u^{(m)}}{v} \right\rceil \\ &\geq [2(v_{n-1}-1)] \\ &= 2(v_{n-1}-1) \end{aligned} \quad (\text{式 (B.7) より})$$

となる。上式より

$$\begin{aligned} 128-2 &\geq 2v_{n-1} \\ \frac{128}{2}-1 &\geq v_{n-1} \end{aligned}$$

すなわち

$$v_{n-1} < \left\lceil \frac{128}{2} \right\rceil = 64$$

となる。

式 (B.2) の仮定のもとで上式が導かれたので、逆に

$$v_{n-1} \geq \left\lceil \frac{128}{2} \right\rceil = 64$$

であれば

$$\hat{q} \leq q_m + 2 \quad (\text{B.1})$$

が成り立つことになる。

## 参 考 文 献

- ( 1 ) D.E.Knuth ( 著 )・中川圭介 ( 訳 )「準数値算法 / 算術演算 : The art of computer programming, 第 4 分冊」, Pp.536、サイエンス社、1986 .
- ( 2 ) D.E.Knuth. Seminumerical algorithms: The art of computer programming, Vol.2, third edition. Pp.762, Addison-Wesley, 1998.
- ( 3 ) 岡本安晴「Delphi プログラミング入門」, Pp.207,CQ 出版株式会社、1997.