

## 15 . 乗除算と固定小数点演算

8ビット表現において、10進数の「10」は2進数では

```
10D=00001010B
```

であり、また、「20」は2進数で

```
20D=00010100B
```

となる。この2つを見比べると「20」の表現は「10」の2進数表現を全体に左に1ビットだけずらした（「シフト (Shift)」という）ものとなっていることがわかる。すなわち、1ビット左シフトすれば「2倍」となるわけである。したがって、2ビット左シフトでは

```
40D=00101000B
```

で「4倍」、3ビット左シフトは「8倍」となる。要するにNビット左シフトは「 $2^N$ 倍」となる。これは2進数が2の累乗で表されている結果である。この時、一番下位のビット (LSB) には「0」が入って来ることに注意されたい。これは、10進数で「10」を10倍するには桁を左に1つずらして「100」とすることと同じである。

このように考えると、逆に右に1ビットシフトすれば「1/2倍」となることも理解できるであろう。すなわち、

```
5D=00000101B
```

となる。さらに、これを右に1ビットシフトすると

```
00000010B=2D
```

となる。この時、「5」の半分は「2.5」であるが、小数部は表現できないので切り捨てとなり、誤差を生じる。また、一番上位ビット (MSB) には「0」を入れることに注意されたい。このようにビットのシフトによって「2の累乗の乗除算」ができる。

いままでは正の数と考えたが、負の数ではどうだろうか。例えば、「-10」は2進数で

```
-10D=11110110B
```

である。これを左1ビットシフトすると

```
-20D=11101100B
```

で2倍となる。負の数においても「 $2^N$ 倍」する時は正の数の場合と同じであり、LSBには「0」をつめている。それでは、「-10」を「1/2倍」するために右に1ビットシフトすると

```
01111011B
```

となる。正の数と同じようにMSBに「0」を入れると、これは「+123」となってしまうことがわかる。「-10」の「1/2倍」は「-5」であり、これは

```
11111011B
```

の2進数表現である。すなわち、負の数の「 $(1/2^N)$ 倍」もまた負の数となるので、MSBは「1」でなければならない。このことは正の数の場合と異なる。このようにして、さらに右1ビットシフトすると

```
11111101B
```

となり、これは「-3」である。実際には「-2.5」であり、正の数の場合とは異なった結果

となる。すなわち、ビットシフト操作で得られる整数値は BASIC における「INT(その値を越えない最大の整数値)」命令実行時と同じ結果である。これらをまとめると

「 $2^N$ 」を得るには「Nビット左シフトしLSBに「0」を入れる  
 「 $1/2^N$ 」を得るには「Nビット右シフトし  
 扱う数値が正の時にはMSBに「0」を入れる  
 扱う数値が正負表現の時はMSBの値は不変とする

ということである。これらの操作を行うアセンブラ命令が8章で述べた「算術シフト命令 (Arithmetic Shift)」であり、「SAL(SHL)およびSAR」命令である。

「2の累乗の乗除算」はシフト操作によって簡単に結果を得ることができ、「乗除算命令 (MUL, DIV 等)」を行うよりはるかに速く実行できるが、乗算(左シフト)の場合には扱う数値の大きさに注意が必要である。

2の累乗計算の場合には前述したような方法で求めることができるが、2の累乗でない乗算を行う場合には「MUL(IMUL)」命令を使用するか、あるいは2の累乗のシフト操作で構成したプログラム、例えば、「AX」の5倍の値を求める時は

$$AX \times 2^2 + AX$$

などのようにしなければならない。後述するDSPでは「乗算命令」は1マシンサイクル(CPUが命令を実行するためのクロック周波数の1周期)で実行できるので、「乗算」に対してはそのまま「乗算命令」を利用しても演算時間の問題はない。ところが、DSPの命令には「除算」がないために除算用のプログラムを組まなければならない(通常30ステップ程度)。したがって、乗算命令を利用して除算するプログラムテクニックによって実行時間の短縮を図る必要があり、ここでこれを説明する。DSP命令は後述するので、CPUで行う場合を述べるが、この方法は実際の制御プログラムで頻繁に使用しているので、必ずマスターされたい。

「5」で除することは「0.2」を乗じることである。「0.2」は「2を乗じて10で除する(あるいは4を乗じて20で除するなど考えてもよく、いくらでもそれらの値は考えられるが)」であるが、これでは乗数も除数も2の累乗ではないために、乗除算命令を使用しなくてはならず、意味がない。そこで、除数を2の累乗に選び、対応する乗数を求める。いま、除数を「 $64 (=2^6)$ 」に選ぶと乗数は「 $64/5=12.8$ 」となるので、この値に近い「13」とする。この乗数を用いて乗算を行うと、例えば「AL=100」に対して

$$\begin{aligned} AX &= 100 \times 13 \\ &= 1300D \\ &= 0000 \ 0101 \ 0001 \ 0100B \end{aligned}$$

という結果が「AX」レジスタに求まる(バイナリ表現はわかり易いように4ビットづつ区切って記述する)。そこで、これを「64」で除するために「AX」を6ビット算術右シフトを行うと

AX = 0000 0000 0001 0100 B

となり、これは「20」である。すなわち、シフトを行う前の「AX」の内容を

AX = 0000 0101 00.01 0100 B

のように、あたかも「6ビットの位置に小数点がある」と考えた時の整数部が答えとして得られている。このような考え方が「固定小数点演算」と呼ばれる演算法の基本である。小数点の右1ビット目は「 $2^{-1}$ 」すなわち「0.5」を表し、右2ビットは「 $2^{-2}$ 」すなわち「0.25」、Nビット目は「 $2^{-N}$ 」を表している。したがって、上記「AX」は「20.3125」を表現している。この値は当然「 $100 \times 13 / 64 = 20.3125$ 」に等しくなる。

さらに、除数の値を大きくして（除数を大きくすると演算誤差が小さくなる）「256 (=  $2^8$ )」とすると、乗数は「 $256 / 5 = 51.2$ 」となるので「51」として、乗算すると

AX = 100 × 51  
 = 5100 D  
 = 0001 0011 1110 1100 B

となる。この時の小数点の位置は

AX = 0001 0011.1110 1100 B

である。ここで、8ビット算術右シフトを行うが、これは「AH」レジスタそのものであり、シフトを行わなくても結果が得られている。ただし、この場合には答えが「19」となっているが、これは乗数「51.2」を「51」と小さい値を用いたためである。小数部を含めると「19.921875」であり、先の乗数を「64」とした時の「20.3125」より「20」に近い値である。正確な値に近い整数値を得るには得られた16ビットデータを「四捨五入」する必要があり、これは固定小数点データに「0.5」を加えて得られる。「0.5」は8ビット小数点表現では

0000 0000.1000 0000 B

であり、これを先のデータに加算すれば

AX = 0001 0100.0110 1100 B

となり、整数部は「20」が得られる（この時のAXの値は20.421875を表している）。実際の制御プログラムで扱っているデータは小数部桁が大きいので四捨五入が行っていないが、より正確な結果を得るにはこの操作が必要である。

このように除数を8ビットでは「256」、16ビットでは「65536」に選べば、結果はそれぞれ8ビットでは「AH」レジスタ、16ビットでは「DX」レジスタに得られる。これは例えば8ビットでは「1」を乗じる場合に「256」を乗じて上位8ビットを取るということであり、単に「1」を「256」に対応させて考えるだけである。扱う数値が正負表現のデータでも同じである。

このような方法は通常「1以上の値で除する、すなわち1以下の値を乗じる」場合には問題がない。しかしながら、時によってはあるデータを「1.3」したいということがあがるが、この時には扱う数値の大きさを考慮していないとオーバーフローする場合があるので、注意を要する。「1.3」倍の時には被乗数8ビットでは乗数が「333」となり、8ビットで表現できないので16ビットに符号拡張して同様に計算すればよいが、被乗数16ビットで

は「85197」となり、計算不可能となる。このような時には次のように演算を工夫する。すなわち、「1.3」倍を

$$1.3 = 1 + 0.3$$

とし、小数部を上述した方法で計算し、その結果に被乗数を加える（1倍は被乗数そのものである）。これは

$$AX \times 1.3 = AX \times 1 + AX \times 0.3$$

である。ただし、小数部の演算結果は上位（この場合は「DX」レジスタ）であるので、「 $AX \times 1$ 」の部分は上位に加算する。このプログラミング例は次のようである（被乗数は「AX」レジスタ、結果は「DX」レジスタ）。

MOV	BX,AX	;AXの内容をBXに保存
MOV	CX,19661	;乗数(0.3×65536)
MUL	CX	
ADD	AX,8000H	:小数部に0.5加算(四捨五入)
ADC	DX,BX	上位にAX×1を加算

扱う数値が正の数だけであれば全ての計算はこれでよいが、正負表現の時には不都合が起こる場合がある。それは乗数の小数部が「0.5」を越える時である（0.5の時には被乗数を1ビット右シフトを行えばよい）。この場合には16ビットデータでは乗数が「32767」より大きくなり、正負表現ではこれを越えると負の値となる。例えば、「1.8」の場合には乗数は「52429」となり、これは上記と同じ計算では「-0.2」になるので、結果は被乗数の「-0.8」倍となる。そこで、このように小数部が「0.5」を越える時には

$$1.8 = 2 - 0.2$$

のように計算し、常に小数部が「0.5」を越えないように行う。「0.2」は「13107 (=65536-52429)」であるから、小数部演算に対して乗数に「13107」を用いると2倍した整数部からの減算、乗数に「52429」を用いると2倍した整数部への加算となる。

これらの演算法は既知の定数で除算する場合に適用できるが、未知数での除算は「DIV(IDIV)」命令を使用しなければならない。

コンピュータで扱う数値は整数であるが、これを整数とみなすか実数とみなすかはプログラマの考え方の違いであり、それによってプログラミングが異なる。整数値とする時、小数点はLSBの右側、すなわち16ビットでは

XXXX	XXXX	XXXX	XXXX	.B
------	------	------	------	----

であると考え。また、前述したように8ビット目に小数点があると考えた場合には

XXXX	XXXX	.XXXX	XXXX	B
------	------	-------	------	---

いずれの表現でも「固定小数点表示 (Fixed-Point Representation)」であるが、後者の場合には小数部も表現されているので、これを特に「分数固定小数点表示 (Fractional Fixed-Point Representation)」という。これらの表現を使用した演算が「固定小数点演算 (Fixed-Point Arithmetic)」であり、制御プログラムではこれを使用している。

一方、実数を仮数部と指数部に分けて表示した場合を「浮動小数点表示 (Floating-Point Representation)」という。これは例えば、16ビットの内、12ビットを仮数部、4ビットを指数部に分けて表示し、

$$(\text{仮数部}) \times 2^{(\text{指数部})}$$

として表す。この表現による演算は行っていないので、詳細はマニュアル等を参考にされたい。

固定小数点について、実際に数値を使って説明するが、わかり易いように電圧、電流と抵抗を使ってオームの法則を計算する。固定小数点の場合、まず、それぞれの値に対して小数点の位置を何処に置くか (小数部のビット数をいくりにするか) を決めなければならない。これにより、表現できる値の大きさと分解能が決まる。ここでは、各値を16ビットで表現することにし、電流を小数部8ビットとする。すなわち

$$\begin{array}{cc} (\text{整数部}) & (\text{小数部}) \\ \text{XXXX XXXX} & \text{.XXXX XXXXB} \end{array}$$

正負表現のため、MSBは符号を表すので、このデータの最大値は

$$0111\ 1111.1111\ 1111\text{B} = 127.99609375\ (\text{A})$$

であり、単位は一応「A」とする。また、最小値は

$$1000\ 0000.0000\ 0000\text{B} = -128\ (\text{A})$$

となり、この範囲の値を表現できることになる。この時の分解能はLSBが「 $2^{-8}$ 」を表すので

$$1\ \text{LSB} = 0.00390625\ (\text{A})$$

となり、この精度で表現可能である。これは当然のことながら

$$128 - 0.00390625 = 127.9960375\ (\text{A})$$

である。実際に実験で流れる電流値は高々20A程度であるから、電流の表現は小数部8ビットで十分である。小数部ビットを大きくすれば精度はよくなるが、それほど精度を上げる必要はない。

抵抗値も電流と同じ精度としておく (単位は「 $\Omega$ 」)。

電圧は実際には300V以上の値を取り得るので、小数部8ビットでは表現できないため、精度を落として小数部4ビットとすると、その最大値は

$$0111\ 1111\ 1111.1111\text{B} = 2047.9375\ (\text{V})$$

であり、単位は一応「V」とする。また、最小値は

$$1000\ 0000\ 0000.0000\text{B} = -2048\ (\text{V})$$

となり、分解能は

$$1\ \text{LSB} = 0.0625\ (\text{V})$$

である。以上のように小数部ビット数を決める。

では、電流  $I=5.5(\text{A})$ 、抵抗  $R=6.1(\Omega)$  の場合の電圧  $V$  を求める。実際の計算では、得られる電圧は

$$V = IR = 5.5 \times 6.1 = 33.55\ (\text{V})$$

である。電流と抵抗の固定小数点表示は

```
電流 I = 1408D = 0580H = 0000 0101.1000 0000 = 5.5 (A)
抵抗 R = 1562D = 061AH = 0000 0110.0001 1010 = 6.1015625 ( )
```

である。ただし、抵抗は 6.1( )に最も近い値を使う。これを「IMUL」命令で積を計算すると、その結果「DX:AX」の 32 ビットデータは

```
00218F00H = 00000000 00100001.10001111 00000000
```

となる。これは、このままの小数点位置であれば

```
33.55859375 (V)
```

を表している。ここで、小数部 8 ビット表現同士の積を行うと、得られる結果は小数部 16 ビット表現となることに注意されたい。電圧は小数部 4 ビットとするので、この 32 ビットから小数部 4 ビットを含む 16 ビットデータを取り出すと(12 ビット右シフトして下位 16 ビットを得るか、あるいは 4 ビット左シフトして上位 16 ビットを得る)

```
00218F00H = 00000000 00100001.10001111 00000000
0000 00100001.1000 = 0218H
```

となる。この値は「33.5(V)」であり、実際値と「0.05(V)」の誤差を生じる。これは演算結果の小数部 8 ビット表現を小数部 4 ビットとしたためであり、分解能による。誤差を小さくするためには小数部のビット数を大きくしなければならない。乗算では小数部 N ビットと M ビットの積は小数部 (N + M) ビット表現の数値となる。

次に、電圧 V=123.7(V)、抵抗 R=8.5( )の場合について電流を計算する。電流の解は

```
I = V/R = 123.7 / 8.5 = 14.55294118 (A)
```

である。電圧と抵抗の固定小数点表示は

```
電圧 V = 1979D = 07BBH = 0000 0111 1011.1011 = 123.6875 (V)
抵抗 R = 2176D = 0880H = 0000 0110.1000 0000 = 8.5 ( )
```

である。ただし、電圧は 123.7(V)に最も近い値を使う。これを計算するために「IDIV」命令を使うが、抵抗は 16 ビット表現なので、電圧を 32 ビット表現に符号拡張して行うと、

```
0000 07BBH / 0880H = 0 (A)
```

となり、正しい結果が得られない。これは小数桁ビットの対応を考慮していないためである。すなわち、乗算とは逆に、除算の場合には被除数が小数部 N ビットで除数が小数部 M ビットでは、商の小数部は (N - M) ビットとなるためである。上の計算では商の小数部ビット数は「-4」となっている。したがって、小数部 8 ビットの商を得るためには、あらかじめ被乗数を 12 ビット左シフトして小数部 16 ビット表現としておく必要がある。

すなわち、電圧は 32 ビット表現で

```
0000 0111 1011.1011
0000 0000 0111 1011.1011 0000 0000 0000 = 007B B000H
```

とする。これで演算を行うと

```
007B B000H / 0880H = 0E8BH
```

であり、これは小数部 8 ビット表現となっているので

```
0000 1110.1000 1101 = 14.55078125 (A)
```

となる。結果の誤差は電圧の表現誤差によるものである。除算の場合には商として得る値

の小数部ビット数を考慮して演算する必要がある。

固定小数点演算においては演算誤差すなわち有効桁数に注意しなければならない。

上述の演算をもう少しわかりやすく記述すると以下のようなものである。電流と抵抗は小数部 8 ビットとしており、これらの 16 ビット値は単に

「1」を「 $256 = 2^8$ 」に対応（実際値の 256 倍）

させ、電圧は

「1」を「 $16 = 2^4$ 」に対応（実際値の 16 倍）

させている。乗算では

$$\begin{aligned} I &= 5.5 \text{ (A)} = 5.5 \times 256 = 1408 \\ R &= 6.1 \text{ ( )} = 6.1 \times 256 = 1561.6 \quad 1562 \end{aligned}$$

とし

$$\begin{aligned} IR &= (1408/256) \times (1562/256) = (1408 \times 1562) / (256 \times 256) \\ &= 2199296 / (256 \times 256) \\ &= 00218F00H / 65536 \end{aligned}$$

であり、結果は実際値の「65536」倍となる。12 ビット右シフトは「 $1/2^{12}$ 」倍すなわち「 $1/4096$ 」倍することであるから

$$V = 2199296 / 4096 = 536.9375 \quad 536 = 0218H$$

となり、実際値の 16 倍の結果が得られる。

一方、除算では

$$\begin{aligned} V &= 123.7 \text{ (V)} = 123.7 \times 16 = 1979.2 \quad 1979 \\ R &= 8.5 \text{ ( )} = 8.5 \times 256 = 2176 \end{aligned}$$

とし、電圧は 12 ビット左算術シフトするので 4096 倍とし

$$V = 1979 \times 4096$$

となり。この時点で電圧は実際値の「 $2^{4+12}$ 」倍となっている。これを 256 倍の抵抗で除するので

$$V/R = (1979 \times 4096) / 2176 = 3725.176471 \quad 3725 = 0E8DH$$

小数部ビット数で表現すれば、乗算では

$$(I \times 2^8) \times (R \times 2^8) = IR \times 2^{16}$$

これを 12 ビット右シフトして

$$IR \times 2^{16-12} = IR \times 2^4$$

となる。除算では、まず電圧を 12 ビット左シフトするので

$$V \times 2^{4+12} = V \times 2^{16}$$

であり、

$$I = (V \times 2^{16}) / (R \times 2^8) = V/R \times 2^{16-8} = V/R \times 2^8$$

となる。

このような演算において、乗算ではオーバーフローは起こらないが、除算では商が 16 ビットに納まるよう、扱う数値の大きさに注意しなければならない。