

ε -pTeX の浮動小数点演算の超簡易説明書

北川 弘典

2008 年 1 月

本文章では、 ε -pTeX で実装されている浮動小数点演算について概説する。詳細な実装方法についてはソースや../ks2/resume.* 他に譲ることにして、ここでは簡単に使い方のみ述べる。

0.1 準備

本文章でいう浮動小数点数とは、よくあるように $-235.673578432E-534$ のように、符号と 10 進小数からなる仮数部に、必要に応じて E or e で始まる指数部が続いたものである。符号は複数あってもよく、そのときは全部掛け算したものが最終的な符号となる。小数点は日米などで使われるピリオドも、欧州大陸で使われるコンマも許容される。

浮動小数点数に纏わるエラーは、オーバーフロー (Floating arithmetic overflow) と例外 (Floating arithmetic exception) の 2 種類のみである。エラーメッセージを出して一旦停止するが、無限大値や NaN 値をそれぞれ代入して続けることができる。

なお、浮動小数点演算を使う場合は、 ε -TeX でいうところの extended mode で処理しなければならない。README.txt にそってフォーマットファイルを作り、それを使った場合は、extended mode は最初から on になっている。

0.2 初期化

浮動小数点演算を行うときには、計算に使用する一時領域や定数 π などを確保する必要がある。これは自動では行われない。

<code>\fpinit</code>	前段落に書いた一時領域その他を確保する。
<code>\fpdest</code>	前段落に書いた一時領域その他を解放する。

初期化を忘れて演算を行おうとしても、特に TeX 側でエラーチェックは行わない。そのため、セグメンテーション違反のようなエラーを引き起こす危険性がある。解放忘れや二重確保は単にメモリを無駄にするだけ*1で、それ以外に問題はない。

0.3 代入, 型変換, 入出力

<code>\real <real></code>	浮動小数点数を表現する glue (以下, $\langle f\text{-glue} \rangle$ と称する) を返す。
---------------------------------	--

*1 それも、TeX の枠内において。

<code>\fpfrac</code> $\langle f\text{-glue} \rangle$	引数の仮数部を返す .
<code>\fpexpr</code> $\langle f\text{-glue} \rangle$	引数の指数部を返す .
<code>\fptoint</code> $\langle f\text{-glue} \rangle$	引数を整数に変換したものを返す . 範囲内に収まらないときは , <code>Number too big</code> エラーを返す . なお , 引数が整数でないときは , 0 に近い方に丸められる .
<code>\fptodim</code> $\langle f\text{-glue} \rangle$	引数を dimension に , 1.0 がちょうど 1pt になるように変換したものを返す . 範囲内に収まらないときは , <code>Dimension too large</code> エラーを返す . なお , 引数が 1/65536 (1sp に対応) でないときは , 同じように 0 に近い方に丸められる .

0.4 四則演算等

<code>\fpadd</code> $\langle f\text{-reg} \rangle$ $\langle real \rangle$	$\langle f\text{-glue} \rangle$ が格納された skip レジスタ (以下 , $\langle f\text{-reg} \rangle$ と称する) と浮動小数点数を引数にとり , 2 つの浮動小数点数の和を計算して , $\langle f\text{-glue} \rangle$ に上書きする .
<code>\fpsub</code> $\langle f\text{-reg} \rangle$ $\langle real \rangle$	同様に差を計算して , $\langle f\text{-glue} \rangle$ に上書きする .
<code>\fpmul</code> $\langle f\text{-reg} \rangle$ $\langle real \rangle$	同様に積を計算して , $\langle f\text{-glue} \rangle$ に上書きする .
<code>\fpdiv</code> $\langle f\text{-reg} \rangle$ $\langle real \rangle$	同様に商を計算して , $\langle f\text{-glue} \rangle$ に上書きする .
<code>\fppow</code> $\langle f\text{-reg} \rangle$ $\langle real \rangle$	同様に累乗を計算して , $\langle f\text{-glue} \rangle$ に上書きする . $x^y = \exp(y \log x)$ で計算するため , 第 1 引数が負ではエラーが生じる .
<code>\fppowi</code> $\langle f\text{-reg} \rangle$ $\langle number \rangle$	<code>\fppow</code> と同様に累乗を計算するが , 第 2 引数 , つまり指数部分は整数に限られる . その代わりに , 第 1 引数は負でもかまわない .

0.5 単項演算

以下は単項演算で , $\langle f\text{-reg} \rangle$ を 1 つとり , 演算をし , 結果をその $\langle f\text{-reg} \rangle$ に上書きする . よって , 以下はコマンドの引数も省略し , 演算内容しか書かない . 引数の内容は便宜的に x で表す .

<code>\fpneg</code>	-1 倍を計算する .
<code>\fpsqr</code>	平方根 \sqrt{x} を計算する .
<code>\fpexp</code>	指数関数 $\exp x$ を計算する .
<code>\fplog</code>	対数関数 $\log x$ を計算する .
<code>\fpabs</code>	絶対値を計算する .
<code>\fpceil</code>	天井関数 $\lceil x \rceil$, つまり x を越えない最小の整数を計算する .
<code>\fpfloor</code>	床関数 $\lfloor x \rfloor$, つまり x 以下の最大の整数を計算する .
<code>\fpsin</code> , ... , <code>\fptan</code>	それぞれ三角関数 $\sin x$, $\cos x$, $\tan x$ を計算する .
<code>\fpsinh</code> , ... , <code>\fptanh</code>	それぞれ双曲線関数 $\sinh x$, $\cosh x$, $\tanh x$ を計算する .
<code>\fpasin</code> , ... , <code>\fpatan</code>	それぞれ逆三角関数 $\arcsin x$, $\arccos x$, $\arctan x$ を計算する . 結果は $\arcsin x$, $\arctan x \in [-\pi/2, \pi/2]$, $\arccos x \in [0, \pi]$ である (主値をとって) .
<code>\fpasinh</code> , ... , <code>\fpatanh</code>	それぞれ逆双曲線関数 $\operatorname{arcsinh} x$, $\operatorname{arccosh} x$, $\operatorname{arctanh} x$ を計算する . 結果は $\operatorname{arcsinh} x$, $\operatorname{arctanh} x \in \mathbf{R}$, $\operatorname{arccosh} x \in [0, \infty]$ の範囲に収まる .

0.6 数値積分によるサンプル

本節では, $f(x) = 1/(x+3)$ を $[-1, 1]$ で, 区間を 40 等分に分割して台形則, 中点則, Simpson 法による数値積分を行う. 当然ながら真値は $\log 2 \simeq 6.93147180559945309417 \times 10^{-1}$ である.

たたき台として, 以下の Fortran 90 のプログラムを使用した. これは 2007 年度夏学期の東京大学理学部数学科の講義「計算数理 I」で僕が提出したレポートの中にあったプログラムを簡略化したものである.

```
PROGRAM main
  IMPLICIT REAL*8 (a-h,o-z)
  a=-1d0; b=1d0; n=40; d=0d0; u=0d0
  DO i=0,n-1
    u=u+1d0/(3d0+a+(b-a)*i/n)
    d=d+1d0/(3d0+a+(b-a)*(i+5d-1)/n)
  END DO
  u=u+5d-1*1d0/(3d0+b)-5d-1*1d0/(3d0+a)
  WRITE(*,*) 'DAIKEI: ', u*(b-a)/n
  WRITE(*,*) 'CHUTEN: ', d*(b-a)/n
  WRITE(*,*) 'SIMPSON: ', (u+2d0*d)*(b-a)/n/3d0
  END
```

これの実行結果は以下である.

```
[h7k doc]$ gfortran -o ks1 ks1.f90
[h7k doc]$ ./ks1
DAIKEI:    0.693186240009141
CHUTEN:    0.693127651979310
SIMPSON:   0.693147181322587
```

これを ϵ -p_{TEX} の浮動小数点演算で書き直して計算させたところ, 以下の結果になった:

台形則での計算結果:	$6.93186240009140538665 \times 10^{-1}$
中点則での計算結果:	$6.9312765197931015099 \times 10^{-1}$
Simpson 則での計算結果:	$6.93147181322586946883 \times 10^{-1}$
真値:	$6.93147180559945309417 \times 10^{-1}$

本文書のソースを示す。ε-TeX の \numexpr 相当の機能がまだ準備されていないので、ソースは無残な姿である。

```
1  %!epllatex fp.tex
   \documentclass[a4j,papersize]{jsarticle}
   \def\epTeX{\varepsilon-\pTeX}\def\etex{\varepsilon-\TeX}
   \def<#1>{\langle\hbox{\it #1}\rangle$}
5  \def\.#1{{\tt\char'134 #1}}
   \def\listx{\def\makelabel{\selectfont } \def\@{\hfill}
   \labelwidth=14zw\labelsep1zw\itemindent11zw\leftmargin=4zw}
   \def\arcsinh{\mathop{\rm arcsinh}}
   \def\arccosh{\mathop{\rm arccosh}}
10 \def\arctanh{\mathop{\rm arctanh}}
   \fpinit % 浮動小数点演算を使用するため
   \usepackage{moreverb}
   \title{\epTeX の浮動小数点演算の超簡易説明書}
   \author{北川 弘典}
15 \date{2008 年 1 月}
   \begin{document}

   \maketitle
   本文章では、\epTeX で実装されている浮動小数点演算について概説する。
20 詳細な実装方法についてはソースや{\tt ../ks2/resume.*}, 他に譲ることにして、
   ここでは簡単に使い方のみ述べる。

   \subsection{準備}
   本文章でいう{\gt 浮動小数点数}とは、よくあるように
25 \.-235.673578432E-534}のように、符号と 10 進小数からなる仮数部に、必要に
   応じて{\tt E} or {\tt e}で始まる指数部が続いたものである。符号は複数あってもよく、
   そのときは全部掛け算したものが最終的な符号となる。小数点は日米などで使わ
   れるピリオドも、欧州大陸で使われるコンマも許容される。

30 浮動小数点数に纏わるエラーは、オーバーフロー ({\tt Floating arithmetic
   overflow}) と例外 ({\tt Floating arithmetic exception}) の 2 種類のみである。エ
   ラーメッセージを出して一旦停止するが、無限大値や NaN 値をそれぞれ代入して
   続けることができる。

35 なお、浮動小数点演算を使う場合は、\etex でいうところの extended mode で処
   理しなければならない。{\tt README.txt}にそってフォーマットファイルを作り、
   それを使った場合は、extended mode は最初から on になっている。

   \subsection{初期化}
40 浮動小数点演算を行うときには、計算に使用する一時領域や定数$\pi$などを確保
   する必要がある。これは自動では行われない。
   \begin{list}{}{\listx}
   \item[.\{fpinit}\@] 前段落に書いた一時領域その他を確保する。
   \item[.\{fpdest}\@] 前段落に書いた一時領域その他を解放する。
45 \end{list}
   初期化を忘れて演算を行おうとしても、特に\TeX 側でエラーチェックは行わない。
```

そのため、セグメンテーション違反のようなエラーを引き起こす危険性がある。
解放忘れや二重確保は単にメモリを無駄にするだけ\footnote{それも、\TeX の枠
内において。}で、それ以外に問題はない。

50

```
\subsection{代入, 型変換, 入出力}
\begin{list}{}{\listx}
\item[\.{freal}\ \<real>\@] 浮動小数点数を表現する glue (以下, \<f-glue>と称
する) を返す .
```

55

```
\item[\.{fpfrac}\ \<f-glue>\@] 引数の仮数部を返す .
\item[\.{fpexpr}\ \<f-glue>\@] 引数の指数部を返す .
\item[\.{fpint}\ \<f-glue>\@] 引数を整数に変換したものを返す . 範囲内に
収まらないときは, {\tt Number too big}エラーを返す . なお, 引数が整数でな
いときは, 0 に近い方に丸められる .
```

60

```
\item[\.{fpdodim}\ \<f-glue>\@] 引数を dimension に, $1.0$ がちょうど
$1\,$pt になるように変換したものを返す . 範囲内に収まらないときは, {\tt
Dimension too large}エラーを返す . なお, 引数が $1/65536$ ( $1\,$sp に対応) で
ないときは, 同じように 0 に近い方に丸められる .
\end{list}$ 
```

65

```
\subsection{四則演算等}
\begin{list}{}{\listx}
\item[\.{fpadd}\ \<f-reg>\ \<real>\@]
\<f-glue> が格納された skip レジスタ (以下, \<f-reg> と称する) と浮動小数点
数を取引数にとり, 2 つの浮動小数点数の和を計算して, \<f-glue> に上書きする .
```

70

```
\item[\.{fpsub}\ \<f-reg>\ \<real>\@]
同様に差を計算して, \<f-glue> に上書きする .
\item[\.{fpmul}\ \<f-reg>\ \<real>\@]
同様に積を計算して, \<f-glue> に上書きする .
```

75

```
\item[\.{fpdiv}\ \<f-reg>\ \<real>\@]
同様に商を計算して, \<f-glue> に上書きする .
\item[\.{fppow}\ \<f-reg>\ \<real>\@] 同様に累乗を計算して, \<f-glue> に
上書きする .  $x^y = \exp(y \log x)$  で計算するため, 第 1 引数が負ではエラーが生じる .
```

80

```
\item[\.{fppow}\ \<f-reg>\ \<number>\@] \<f-reg> と同様に累乗を計算する
が, 第 2 引数, つまり指数部分は整数に限られる .
その代わりに, 第 1 引数は負でもかまわない .
\end{list}
```

```
\subsection{単項演算}
```

85

以下は単項演算で, \<f-reg> を 1 つとり, 演算をし, 結果をその \<f-reg> に上
書きする . よって, 以下はコマンドの引数も省略し, 演算内容しか書かない . 引
数の内容は便宜的に x で表す .

```
\begin{list}{}{\listx}
```

90

```
\item[\.{fpneg}\@]  $-1$  倍を計算する .
\item[\.{fpsqr}\@] 平方根  $\sqrt{x}$  を計算する .
\item[\.{fpexp}\@] 指数関数  $\exp x$  を計算する .
\item[\.{fplog}\@] 対数関数  $\log x$  を計算する .
\item[\.{fpabs}\@] 絶対値を計算する .
```

95

```
\item[\.{fpceil}\@] 天井関数  $\lceil x \rceil$ , つまり  $x$  を越えない最小の
```

整数を計算する .

\item[\.{fplfloor}\@] 床関数 $\lfloor x \rfloor$, つまり x 以下の最大の整数を計算する .

\item[\.{fpsin}, \dots, \.{fptan}\@]

100 それぞれ三角関数 $\sin x$, $\cos x$, $\tan x$ を計算する .

\item[\.{fpsinh}, \dots, \.{fptanh}\@]

それぞれ双曲線関数 $\sinh x$, $\cosh x$, $\tanh x$ を計算する .

\item[\.{fpasin}, \dots, \.{fpatan}\@] それぞれ逆三角関数 $\arcsin x$, $\arccos x$, $\arctan x$ を計算する . \ 結果は $\arcsin x$, \arctan

105 $x \in [-\pi/2, \pi/2]$, $\arccos x \in [0, \pi]$ である (主値をとって) .

\item[\.{fpasinh}, \dots, \.{fpatanh}\@] それぞれ逆双曲線関数 $\operatorname{arcsinh} x$, $\operatorname{arccosh} x$, $\operatorname{arctanh} x$ を計算する . \ 結果は $\operatorname{arcsinh} x$, $\operatorname{arctanh} x \in \mathbb{R}$, $\operatorname{arccos} x \in [0, \infty)$ の範囲に収まる .

\end{list}

110

\subsection{数値積分によるサンプル}

本節では, $f(x)=1/(x+3)$ を $[-1, 1]$ で, 区間を40等分に分割して台形則, 中点則, Simpson法による数値積分を行う . 当然ながら真値は

\skip300=\real2\fplog\skip300\log2\simeq\fpfrac\skip300\times

115 $10^{\text{fpexpr}\text{skip300}}$

である .

たたき台として, 以下のFortran 90のプログラムを使用した . これは2007年度夏学期の東京大学理学部数学科の講義「計算数理 I」で僕が提出したレポートの中

120 にあったプログラムを簡略化したものである .

\begin{verbatim}

PROGRAM main

IMPLICIT REAL*8 (a-h,o-z)

a=-1d0; b=1d0; n=40; d=0d0; u=0d0

125 DO i=0,n-1

u=u+1d0/(3d0+a+(b-a)*i/n)

d=d+1d0/(3d0+a+(b-a)*(i+5d-1)/n)

END DO

u=u+5d-1*1d0/(3d0+b)-5d-1*1d0/(3d0+a)

130 WRITE(*,*) 'DAIKEI: ', u*(b-a)/n

WRITE(*,*) 'CHUTEN: ', d*(b-a)/n

WRITE(*,*) 'SIMPSON: ', (u+2d0*d)*(b-a)/n/3d0

END

\end{verbatim}

135 これの実行結果は以下である .

\begin{verbatim}

[h7k doc]\$ gfortran -o ks1 ks1.f90

[h7k doc]\$./ks1

DAIKEI: 0.693186240009141

140 CHUTEN: 0.693127651979310

SIMPSON: 0.693147181322587

\end{verbatim}

145 %以降にプログラムに入る．とりあえずは上のものを逐語訳する方向でいこう．
これを\epTeX の浮動小数点演算で書き直して計算させたところ，以下の結果に
なった：

```

\par\vskip0.5\baselineskip\par
150 \newskip\nia\newskip\nib\newskip\nid
\newskip\niu\newcount\nin\newcount\nii
%
\nia=\real-1 \nib=\real1 \nin=40 \nid=\real0 \niu=\nid % line 3
%
155 \nii=0
\loop \ifnum\nii<\nin\relax
\skip300=\real3 \fpadd\skip300\nia % \skip300 = 3+a
\skip301=\real\nii \fpdiv\skip301\nin % \skip301=i/n
\skip302=\nib \fsub\skip302\nia \fpmul\skip302\skip301 % \skip302=(b-a)*i/n
160 \skip301=\skip300 \fpadd\skip301\skip302 % \skip301=3+a+(b-a)*i/n
\fpowi\skip301 by -1 \fpadd\niu\skip301 % line 5
\skip301=\real\nii \fpadd\skip301by0.5 \fpdiv\skip301\nin
\skip302=\nib \fsub\skip302\nia \fpmul\skip302\skip301
\skip301=\skip300 \fpadd\skip301\skip302
165 \fpowi\skip301 by -1 \fpadd\nid\skip301 % line 6
\advance\nii by1
\repeat
%
\skip300=\real0.5 \skip301=\real3 \fpadd\skip301\nib
170 \fpdiv\skip300\skip301
\skip301=\real0.5 \skip302=\real3 \fpadd\skip302\nia
\fpdiv\skip301\skip302
\fsub\skip300\skip301 \fpadd\niu\skip300 % line 8
%
175 \skip300=\nib\fsub\skip300\nia\fpdiv\skip300by\nin
\fpmul\niu\skip300 \fpmul\nid\skip300 % 先に (b-a)/n で掛けておく
%
\noindent
\leavevmode\hbox to 13zw{台形則での計算結果:\hss}%
180 $\fpfrac\niu\times 10^{\fpexpr\niu}$\
%
\leavevmode\hbox to 13zw{中点則での計算結果:\hss}%
$\fpfrac\nid\times 10^{\fpexpr\nid}$\
%
185 \leavevmode\hbox to 13zw{Simpson 則での計算結果:\hss}%
\skip300=\niu\fpadd\skip300\nid\fpadd\skip300\nid\fpdiv\skip300by3
$\fpfrac\skip300\times 10^{\fpexpr\skip300}$\
%
\leavevmode\hbox to 13zw{真値:\hss}%
190 \skip300=\real2\fplog\skip300
$\fpfrac\skip300\times 10^{\fpexpr\skip300}$

\newpage

```

本文書のソースを示す。`\eTeX` の`\verb+\numexpr+` 相当の機能がまだ準備されていないので、ソースは無残な姿である。

```
\small  
\listinginput[5]{1}{fp.tex}  
\end{document}
```