

数理解析学特論III講義資料

高安亮紀 (筑波大学)*

2018年11月16日

1. 浮動小数点数と誤差

計算機で扱える数「整数」と「浮動小数点数」について解説し、「丸め誤差」などの数値計算の誤差と、それを制御する「区間演算」を紹介する。

そもそも数値計算とは？

- 数値計算（数値解析）とは…数学の問題を有限桁の浮動小数点数を使って「数値的」に解く。数式処理とは違う
- 数値計算でできること，できないこと
- 数値計算は誤差との戦い
- 整数（integer）と浮動小数点数（binary32/binary64）

1.1. 整数型

Matlab上で整数は

型	機械内表現
int8	8bit 整数
int16	16bit 整数
int32	32bit 整数
int64	64bit 整数

のような数値がある。

```
a = int32(10); disp(a)
```

```
10
```

```
b = int32(2147483647)+int32(1); disp(b)
```

```
2147483647
```

```
c = intmax('int32'); disp(c)
```

```
2147483647
```

これからint32の最大値が2147483647であることがわかる。実際 $2^{31}-1 = 2147483647$ である。これはint32整数が次のようなビットパターンで表現されていることによる。

* 〒1305-8573 茨城県つくば市天王台1-1-1
e-mail: takitoshi@risk.tsukuba.ac.jp
web: <http://www.risk.tsukuba.ac.jp/~takitoshi/>

ビットパターン	数値
01111111111111111111111111111111	2147483647
00000000000000000000000000000010	2
0000000000000000000000000000000001	1
0000000000000000000000000000000000	0
11111111111111111111111111111111	-1
11111111111111111111111111111110	-2
10000000000000000000000000000000	-2147483648

このような負の数の表現形式を「2の補数形式」と呼ぶ。32個の各bitが次のような重みをもっていると考えられる。

$$\boxed{-2^{31}} \boxed{2^{30}} \boxed{2^{29}} \dots \boxed{2^1} \boxed{2^0}$$

2の補数形式の場合、 n ビットで $-2^{n-1} \sim 2^{n-1} - 1$ の範囲の数を表現できる。つまり

型	表現範囲
int8	-128～127
int16	-32768～32767
int32	-2147483648～2147483647
int64	-9223372036854775808～9223372036854775807

の範囲の整数が表せる。

```
b = int64(2147483647)+1; disp(b)
```

```
2147483648
```

```
b = int64(9223372036854775807)+1; disp(b)
```

```
9223372036854775807
```

1.2. 浮動小数点数

浮動小数点数 (binary64, binary32) は、「浮動小数点形式」と呼ばれる形式で表現できる。例えば、「1234.5」を「 1.2345×10^3 」のように、小数点の位置を1番左の数値と左から2番目の数値の間に移動（「正規化」と呼ぶ）し、それに指数を掛けた形式で数を表現する。この「1.2345」の部分を「仮数部」、 10^3 の部分（厳密には³）を「指数部」という。

浮動小数点数は仮数部の長さ、指数部の長さ、指数が2, 10, 16など、多様な規格が考えられる。そこで1985年にWilliam Kahanが中心となって

IEEE 754: Standard for Binary Floating-Point Arithmetic
という標準規格が制定された。最近では世に出るハードウェアのほぼ全てがこの規格に従っている。

1.2.1. 倍精度 (binary64)

倍精度は、符号 (±) に1ビット、指数部に11ビット、仮数部に52ビットを使う。全部で64ビット=8バイトである。

1 (符号)	11 (指数部)	52 (仮数部)
--------	----------	----------

- 符号は、0なら正、1なら負
- 指数部は「 $\times 2^{\text{指数}}$ 」の指数の部分に1023を加えたものが11ビット符号無し of 整数の形で格納されている
- 仮数部は、実際の仮数部の先頭の「1」を取り除いた残りが格納されている。仮数部の先頭は必ず1にしてメモリに格納しないことで1ビット分精度を稼いでいる

数値 x は

$$x = \pm 1.d_1d_2 \cdots d_{52} \times 2^m = \pm \left(\frac{1}{2^0} + \frac{d_1}{2^1} + \frac{d_2}{2^2} + \cdots + \frac{d_{52}}{2^{52}} \right) 2_{(10)}^e$$

と書ける ($-1022 \leq e \leq 1023$, $m : e + 1023$ の2進表現)。

例えば、5.25は2進数で書くと

$$101.01_{(2)} = \left(\frac{1}{2^0} + \frac{0}{2^1} + \frac{1}{2^2} + \frac{0}{2^3} + \frac{1}{2^4} \right) \times 2_{(10)}^2$$

であるから、計算機内では

```
0 | 10000000001 | 0101000000000000000000000000000000000000000000000000000
```

のように格納されている。指数部の「1000000001」は、「 $2+1023=1025$ 」を2進数にしたもの。これは正規化数と呼ばれる数の表現。

`x = 5.25;`

```
disp(float2bin(x))
```

```
disp(float2bin(5.24))
```

```
0100000000010101000000000000000000000000000000000000000000000000000000
```

```
0100000000010100111101011100001010001111010111000010100011110110
```

```
binx = float2bin(x);
```

```
sign = binx(1); disp(sign)
```

```
exp = binx(2:12); disp(exp)
```

```
frac = binx(13:end); disp(frac)
```

```
0
```

```
10000000001
```

```
0101000000000000000000000000000000000000000000000000000000000000
```

符号は `sign = 0` で正の数、指数部は `exp-1023=1025-1023=2`、仮数部は

```
frac=(1010000000000000000000000000000000000000000000000000000000000000)_2
```

で2進数表示され、実際の値は1.3125。元の浮動小数点数に戻すと

非正規化数

指数部が $e + 1023 = 0$ かつ仮数部が0でないとき、仮数部の最初の桁を0にして

$$\pm 0.d_1d_2 \cdots d_{52} \times 2^0 = \pm \left(\frac{0}{2^0} + \frac{d_1}{2^1} + \frac{d_2}{2^2} + \cdots + \frac{d_{52}}{2^{52}} \right) 2^{-1022} \quad (10)$$

という数の表現をする。非正規化数は文字通り「正規化していない」数。

```
x = 2^(-1022); disp(x), disp(float2bin(x))
disp(float2bin(x-2^(-1074)))
```

```
2.225073858507201e-308
00000000000100000000000000000000000000000000000000000000000000000000000000000000000000
0000000000000111111111111111111111111111111111111111111111111111111111111111111111111111111
```

漸近アンダーフロー

正規化数の最小数の最終 bit を1だけ減じると

$$\left(\frac{0}{2^0} + \frac{1}{2^1} + \frac{1}{2^2} + \cdots + \frac{1}{2^{52}} \right) 2^{-1022} \quad (10)$$

となり、これを正規化すると

$$\left(\frac{1}{2^0} + \frac{1}{2^1} + \frac{1}{2^2} + \cdots + \frac{0}{2^{52}} \right) 2^{-1023} \quad (10)$$

となって、指数部の下限 $-1022 \leq e$ を超えてしまう。そこで、「 2^{-1022} を下回ったら正規化をやめて指数部を 2^{-1022} に固定して仮数部の最初の桁を0としてみて格納する」とする。これが非正規化数であり、

$$\left(\frac{0}{2^0} + \frac{1}{2^1} + \frac{1}{2^2} + \cdots + \frac{1}{2^{52}} \right) 2^{-1022} \quad (10)$$

$$\left(\frac{0}{2^0} + \frac{1}{2^1} + \frac{1}{2^2} + \cdots + \frac{0}{2^{52}} \right) 2^{-1022} \quad (10)$$

⋮

$$\left(\frac{0}{2^0} + \frac{1}{2^1} + \frac{0}{2^2} + \cdots + \frac{0}{2^{52}} \right) 2^{-1022}$$

$$\left(\frac{0}{2^0} + \frac{0}{2^1} + \frac{1}{2^2} + \cdots + \frac{1}{2^{52}} \right) 2^{-1022} \quad (10)$$

⋮

$$\left(\frac{0}{2^0} + \frac{0}{2^1} + \frac{0}{2^2} + \cdots + \frac{1}{2^{52}} \right) 2^{-1022} = 2^{-1074} \approx 10^{-323.31} \quad (10)$$

のような数が表現できる。ただし、 2^{-1022} と 2^{-1074} の間の数は、本来53ビットあるべき仮数部の長さが52ビット～1ビットまで減ってしまっており、精度が低下していることに注意。

```

x = 2^(-1074); disp(x), disp(float2bin(x))
disp(x/2), disp(float2bin(x/2))

4.940656458412465e-324
000000000000000000000000000000000000000000000000000000000000000000000000000000000000000000000000001
0
000000000000000000000000000000000000000000000000000000000000000000000000000000000000000000000000000

```

1.2.2. 浮動小数点数まとめ

倍精度 (binary64) は

	仮数部が0	仮数部が0でない
$e + 1023 = 0$	± 0	非正規化数
$1 \leq e + 1023 \leq 2046$	正規化数	正規化数
$e + 1023 = 2047$	$\pm \infty$	NaN

単精度 (binary32) は

	仮数部が0	仮数部が0でない
$e + 127 = 0$	± 0	非正規化数
$1 \leq e + 127 \leq 254$	正規化数	正規化数
$e + 127 = 255$	$\pm \infty$	NaN

今後、浮動小数点数全体の集合を \mathbb{F} と表すことにする。特に断りがなければ、浮動小数点数は倍精度浮動小数点数 (64bit) とする。IEEE754 では **binary64** と呼ばれている。

1.3. 丸め誤差

浮動小数点数を使用した演算に混入する誤差について述べる。浮動小数点数同士の演算 (加減乗除など) の結果は、浮動小数点数で表せるとは限らない。例えば、10進数で仮数部3桁の浮動小数点演算を考え、 $2/3$ を計算すると、

$$2.00 \times 10^0 / 3.00 \times 10^0 = 0.66666666\dots \times 10^0$$

となり、仮数部3桁に収まらない。仮数部の4桁目で四捨五入を行うと

$$6.67 \times 10^{-1}$$

となる。このときの計算値と真値との差

$$6.67 \times 10^{-1} - 6.6666666\dots \times 10^{-1} = 3.3333333 \times 10^{-4}$$

が丸め誤差である。

- 最近点への丸め（デフォルト）： \tilde{x} に最も近い浮動小数点数に丸める。もし2点あるならば仮数部の最後のビットが0である浮動小数点数に丸める。
- $+\infty$ 方向への丸め： \tilde{x} 以上の浮動小数点数の中で最も小さい浮動小数点数に丸める。
- $-\infty$ 方向への丸め： \tilde{x} 以下の浮動小数点数の中で最も大きい浮動小数点数に丸める。
- 原点方向への丸め：絶対値が \tilde{x} 以下の浮動小数点数の中で、 \tilde{x} に最も近いものに丸める。

1.5. その他の誤差いろいろ

丸め誤差は実数を浮動小数点数で近似する際の不正確さであった。ここでは浮動小数点数を用いた演算の問題点を紹介する。

1.5.1. 桁落ち

極めて近い数どうしの減算によって、誤差が著しく大きくなってしまう現象。2つの浮動小数点数

$$x = \left(\frac{1}{2^0} + \frac{d_1}{2^1} + \cdots + \frac{d_p}{2^p} + \frac{1}{2^{p+1}} + \cdots + \frac{b_{p+2}}{2^{p+2}} + \cdots + \frac{b_{52}}{2^{52}} \right) 2^e_{(10)},$$

$$y = \left(\frac{1}{2^0} + \frac{d_1}{2^1} + \cdots + \frac{d_p}{2^p} + \frac{0}{2^{p+1}} + \cdots + \frac{c_{p+2}}{2^{p+2}} + \cdots + \frac{c_{52}}{2^{52}} \right) 2^e_{(10)}$$

が $x > y$ とし、仮数部の最初から p ビットが等しいとする。このとき

$$x - y = \left(\frac{1}{2^0} + \frac{b_{p+2} - c_{p+2}}{2^1} + \cdots + \frac{b_{52} - c_{52}}{2^{52-p-1}} \right) 2^{e-p-1}_{(10)}$$

これよりもともと 52 個あった仮数部の情報が、 $52 - p$ 個に減っている。例を挙げよう。 $b > 0$ とし、2次方程式 $x^2 + bx + c = 0$ の解の公式

$$x_1 = \frac{-b + \sqrt{b^2 - 4c}}{2}, \quad x_2 = \frac{-b - \sqrt{b^2 - 4c}}{2}$$

を考える。いまもしも $b^2 \gg c$ となるならば、 b と $\sqrt{b^2 - 4c}$ が近い数になるので、 x_1 の分子の計算で、桁落ちが起こる。

`b = 1e15; c = 1e14;`

`x1 = (-b+sqrt(b^2-4*c))/2; disp(x1), disp(x1^2+b*x1+c)`

`x2 = 2*c/(-b-sqrt(b^2-4*c)); disp(x2), disp(x2^2+b*x2+c)`

`-0.12500000000000000`

`-2.4999999999999998e+13`

`-0.10000000000000000`

`0`

1.5.2. 情報落ち

絶対値の大きさが極端に違う2数の加減算を行った時、小さいほうの数値の下位の桁が失われてしまう現象.

```
disp((3.14159265358979+1e10)-1e10)
disp(3.14159265358979 + (1e10-1e10))
```

```
3.141592025756836
3.141592653589790
```

```
disp(1e48+543.2-1e48-1e36+123.4+1e36)
disp(1e48-1e48-1e36+1e36+543.2+123.4)
disp(1e48-1e36-1e48+1e36+543.2+123.4)
```

```
0
6.666000000000000e+02
-3.923098441916162e+30
```

1.5.3. 打ち切り誤差

無限回行うべき計算を有限回の計算で置き換えることにより生じる誤差. 計算機は有限回の四則演算しかできない. そのため, 無限級数や収束列のような値を求めるためには, 有限項で打ち切った近似値を用いる. その際に誤差が生じる. (例) Taylor展開の打ち切り誤差, Newton法の打ち切り誤差, 数値積分の打ち切り誤差

```
disp(exp_taylor(20))
disp(exp(20))
disp(exp_taylor(-20))
disp(exp(-20))
```

```
function s = exp_taylor(x)
    s = 0; t = 1; i = 1;
    while 1
        s = s + t;
        if abs(t) < abs(s) * 1e-15, break; end
        t = (x / i)*t;
        i = i+1;
    end
end
```

```
4.851651954097903e+08
4.851651954097903e+08
6.147561828914626e-09
2.061153622438558e-09
```

2. 区間演算

区間の表現, 区間演算について述べる. 精度保証付き数値計算の基本的な原理は, 実数値で与えられる真の値の上限と下限を浮動小数点数により浮動小数点演算を用いて包み込むことである.

\mathbb{R} 上の区間を

$$\mathbf{a} := \{x \in \mathbb{R} : \underline{a} \leq x \leq \bar{a}, \underline{a}, \bar{a} \in \mathbb{R}\} = [\underline{a}, \bar{a}]$$

と表し, 区間の全体を \mathbb{IR} とする. この時 \underline{a} を区間の下端, \bar{a} を区間の上端という. さらに区間の

$$\begin{aligned} \text{直径 (diameter)} : \quad d(\mathbf{a}) &= \bar{a} - \underline{a} \in \mathbb{R} \\ \text{半径 (radius)} : \quad \text{rad}(\mathbf{a}) &= \frac{\bar{a} - \underline{a}}{2} \in \mathbb{R} \\ \text{中点 (center)} : \quad \text{mid}(\mathbf{a}) &= \frac{\bar{a} + \underline{a}}{2} \in \mathbb{R} \\ \text{最小絶対値} : \quad \text{mig}(\mathbf{a}) &= \min\{|a| : a \in \mathbf{a}\} \in \mathbb{R} \\ \text{最大絶対値} : \quad \text{mag}(\mathbf{a}) &= \max\{|a| : a \in \mathbf{a}\} \in \mathbb{R} \end{aligned}$$

をそれぞれ表すとする.

```
a = infsup(-1,2); disp(a)
disp(['d(a) = ', num2str(diam(a))])
disp(['rad(a) = ', num2str(rad(a))])
disp(['mid(a) = ', num2str(mid(a))])
disp(['mig(a) = ', num2str(mig(a))])
disp(['mag(a) = ', num2str(mag(a))])

[ -1.000000000000000,  2.000000000000000]
d(a) = 3
rad(a) = 1.5
mid(a) = 0.5
mig(a) = 0
mag(a) = 2
```

2.1. 区間演算 (上端下端型)

区間 $\mathbf{X} = [a, b]$, $\mathbf{Y} = [c, d]$ に対して, 四則演算を定義する:

$$\begin{aligned} \mathbf{X} + \mathbf{Y} &= [a + c, b + d] \\ \mathbf{X} - \mathbf{Y} &= [a - d, b - c] \\ \mathbf{X} \times \mathbf{Y} &= [\min\{ac, ad, bc, bd\}, \max\{ac, ad, bc, bd\}] \\ \mathbf{X} \div \mathbf{Y} &= [\min\{a/c, a/d, b/c, b/d\}, \max\{a/c, a/d, b/c, b/d\}] \end{aligned}$$

\times, \div は場合分けをする.

	$c > 0$	$0 \in \mathbf{Y}$	$d < 0$
$a > 0$	$[ac, bd]$	$[bc, bd]$	$[bc, ad]$
$0 \in \mathbf{X}$	$[ad, bd]$	A	$[bc, ac]$
$b < 0$	$[ad, bc]$	$[ad, ad]$	$[bd, ac]$

ただし $A = [\min\{ad, bc\}, \max\{ad, bc\}]$.

注意 2.1. 区間内全ての要素について演算を行うため無限回の計算が必要のように思えるが、これをまとめて実現するのが区間演算。さらに除算で割る区間 Y に 0 が含まれる ($0 \in Y$) と演算結果は無限大を含むことになる。

```
X = infsup(-2,1);
Y = infsup(1,3);
disp(X+Y)
disp(X-Y)
disp(X*Y)
disp(X/Y)
Z = Y/X; disp(Z)
```

```
[ -1.0000000000000000,  4.0000000000000000]
[ -5.0000000000000000,  0.0000000000000000]
[ -6.0000000000000000,  3.0000000000000000]
[ -2.0000000000000000,  1.0000000000000000]
+/-Inf
```

最後の +/-Inf は区間 $[-\infty, \infty]$ を表している。

包含関係の単調性が成立する。すなわち $X_1 \subseteq X_2$ かつ $Y_1 \subseteq Y_2$ ならば $X_1 \circ Y_1 \subseteq X_2 \circ Y_2$ が成立する。ここで、 $\circ \in \{+, -, \times, \div\}$ とする。

さらに $+$, \times に関しては、交換則と結合則が成立する。

$$X \circ Y = Y \circ X, X \circ (Y \circ Z) = (X \circ Y) \circ Z, \circ \in \{+, \times\}.$$

しかし、加法と乗法の逆元は存在しない。すなわち分配則が成立しない。劣分配則のみ成立する。

$$X(Y + Z) \subseteq XY + XZ.$$

```
X = infsup(-1,1);
Y = infsup(1,2);
Z = infsup(-2,1);
disp(X*(Y+Z))
disp(X*Y+X*Z)
disp(in(X*(Y+Z),X*Y+X*Z))
```

```
[ -3.0000000000000000,  3.0000000000000000]
[ -4.0000000000000000,  4.0000000000000000]
```

2.2. 中心半径型区間演算

上端下端型区間の他に区間を中心と半径で表す型もある。それを中心半径型区間という。中心半径型 (mid-rad 型) 区間は実数 $x \in \mathbb{R}$ に対して、区間の中心を x_c 、半径を x_r としたとき

$$\mathbf{x} = \langle x_c, x_r \rangle = \{x : x_c - x_r \leq x \leq x_c + x_r\}$$

と表される。中心半径型区間と上端下端型区間の間には次のような関係が成り立つ：

$$\mathbf{x} = [\underline{x}, \bar{x}] = [x_c - x_r, x_c + x_r] = \left\langle \frac{x + \bar{x}}{2}, \frac{x - \bar{x}}{2} \right\rangle = \langle x_c, x_r \rangle = \mathbf{x}$$

中心半径型区間 $\mathbf{x} = \langle x_c, x_r \rangle$, $\mathbf{y} = \langle y_c, y_r \rangle$ に対して、四則演算を定義する。

$$\begin{aligned} \mathbf{x} + \mathbf{y} &= \langle x_c + y_c, x_r + y_r \rangle \\ \mathbf{x} - \mathbf{y} &= \langle x_c - y_c, x_r + y_r \rangle \\ \mathbf{x} \times \mathbf{y} &\subseteq \langle x_c y_c, |x_c| y_r + |y_c| x_r + x_r y_r \rangle \\ \mathbf{x} \div \mathbf{y} &= \frac{\mathbf{x} \times \mathbf{y}}{y \bar{y}} = \frac{\mathbf{x} \times \mathbf{y}}{y_c^2 - y_r^2} \end{aligned}$$

ちなみに乗算は $\text{sgn}()$ を括弧内の符号を返す関数とすると

$$\mathbf{x} \times \mathbf{y} = \langle x_c y_c + \delta_1, \delta_2 \rangle$$

と書ける。ここで

$$\begin{aligned} \delta_1 &= \text{sgn}(x_c y_c) \min \{x_r |y_c|, |x_c| y_r, x_r y_r\}, \\ \delta_2 &= \max \{x_r (|y_c| + y_r), (|x_c| + x_r) y_r, x_r |y_c| + |x_c| y_r\} \end{aligned}$$

である。

```
X = midrad(0,1);
Y = midrad(1.5,0.5);
Z = midrad(-0.5,1.5);
format midrad
disp(X+Y)
disp(X-Y)
disp(X*Y)
disp(X/Y)
```

```
< 1.5000000000000000, 1.5000000000000000>
< -1.5000000000000000, 1.5000000000000000>
< 0.0000000000000000, 2.0000000000000000>
< 0.0000000000000000, 1.0000000000000000>
```

2.3. 関数の値域評価

$I \in \mathbb{I}\mathbb{R}$ をある区間とし, $f: D \subset \mathbb{R} \rightarrow \mathbb{R}$ を領域 D で連続な関数とする. このとき関数 f は $\mathbb{I}\mathbb{R}$ 上の関数として拡張できる.

$$f(I) = \{f(x) : x \in I\}$$

この $f(I)$ は関数の値域で, これを厳密に計算することは非線形関数の場合不可能である. よって

$$f(I) \subseteq [a, b]$$

となる区間 $[a, b]$ で関数の値域を包含する. これを f の区間 I における区間拡張 (区間拡張は一意的でないことに注意!) といい, $f_{\square}(I) = [a, b]$ と表す.

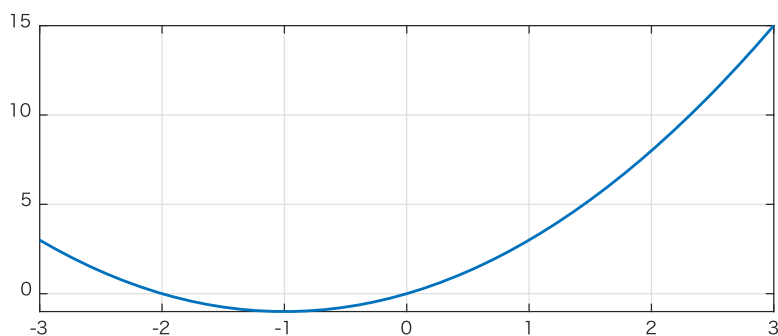
注意 2.2. IEEE 754 標準規格における丸めに従っているのは「四則演算・平方根」だけ. それ以外の関数 (例えば, \sin , \cos などは) 丸めの向きの変更が不可能 (自作するか, 丸めに対応し作成されたものを使うこと). ただし, デフォルトの関数も精度が十分に高く作られているため, 精度保証などを考えない場合はそれで十分.

2.3.1. 例

例えば, 関数 $f(x) = x^2 + 2x$ を考える. 関数に区間を単純に代入すると

$$\begin{aligned}x &\in [0.9, 1.1] \\x^2 &\in [0.81, 1.21] \\2x &\in [1.8, 2.2] \\x^2 + 2x &\in [2.61, 3.41]\end{aligned}$$

```
format infsup
f = @(x) x.^2+2*x;
LW = 'linewidth'; lw = 1.6;
fplot(f, [-3,3], LW, lw)
f(intfsup(0.9,1.1))
```



```
intval ans =
[ 2.609999999999999, 3.410000000000001]
```

```
disp(f(0.9))
disp(f(1.1))
```

2.6100000000000000

3.4100000000000000

$f(x) = x^2 + 2x$ のときは値域評価が区間拡張を利用しても精度が良い。一方で $f(x) = x^2 - 2x$ のときは

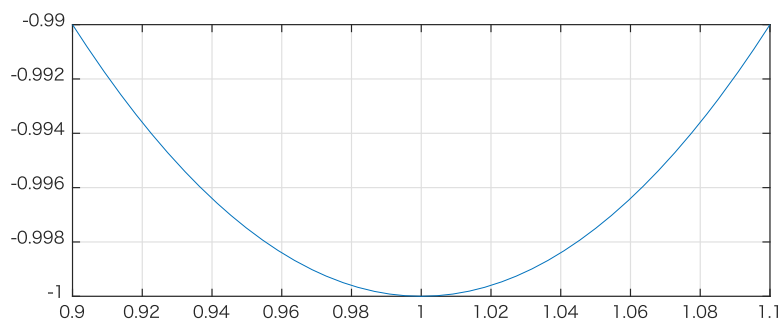
$$x \in [0.9, 1.1]$$

$$x^2 \in [0.81, 1.21]$$

$$2x \in [1.8, 2.2]$$

$$x^2 - 2x \in [-1.39, -0.59] \text{ (幅 } 0.8\text{!)}$$

```
f = @(x) x.^2-2*x;
fplot(f, [0.9, 1.1])
disp(f(0.9))
disp(f(1))
disp(f(1.1))
```



-0.9900000000000000

-1

-0.9900000000000000

これから値域は $[-1, -0.99]$ ，一方で区間演算による区間拡張は $[-1.39, -0.59]$ 。だいぶ過大評価になってしまう。そこで区間幅を改良したい。

アイデア 1. $f(x) = x(x - 2)$;

$$x \in [0.9, 1.1]$$

$$x(x - 2) = [0.9, 1.1] \times [-1.1, -0.9] = [-1.21, -0.81] \text{ (幅 } 0.4\text{)}$$

大分改善された。もう一声！

アイデア 2. $f(x) = (x - 1)^2 - 1$;

$$x \in [0.9, 1.1]$$

$$(x - 1)^2 - 1 = [-0.1, 0.1]^2 - 1$$

$$= [0, 0.01] - 1$$

$$= [-1, -0.99] \text{ (幅 } 0.01\text{)}$$


```

x = infsup(0.9,1.1);
disp(x^2-2*x)
disp(x*(x-2))
disp((x-1)^2-1)

[ -1.390000000000001, -0.589999999999999]
[ -1.210000000000001, -0.809999999999999]
[ -1.000000000000000, -0.989999999999999]

```

2.3.2. 区間幅の抑制のための工夫

1. 区間分割

$$I = I_1 \cup I_2$$

と考えると, $f_{\square}(I_1) \cup f_{\square}(I_2)$ を計算する.

2. 平均値形式

区間拡張 $f_{\square}(I)$ を

$$f_{\square}(I) = f(c) + f'_{\square}(I)(I - c), \quad c = \text{mid}(I)$$

によって得る. ただし $f'_{\square}(I)$ は区間 I における f の1階微分の区間拡張.

2.4. 機械区間演算

区間演算をコンピュータで実現するには \mathbb{R} の代わりに \mathbb{F} を使った区間が必要. そのような区間全体を

$$\mathbb{IF} := \{x \in \mathbb{IR} : \underline{x}, \bar{x} \in \mathbb{F}\}$$

と定義する. IEEE754規格に準拠したシステム上では演算後の丸めの向きを制御することができる. 演算結果が浮動小数点数でない場合, 丸めの向きを制御して計算する. いま $a, b \in \mathbb{F}$ に対して, $\circ \in \{+, -, \times, \div\}$ として

$$\text{fl}_{\nabla}(a \circ b) := \max\{x \in \mathbb{F} : x \leq a \circ b\} \text{ (下向き丸め)}$$

$$\text{fl}_{\Delta}(a \circ b) := \min\{x \in \mathbb{F} : x \geq a \circ b\} \text{ (上向き丸め)}$$

とすると

$$\text{fl}_{\nabla}(a \circ b) \leq a \circ b \leq \text{fl}_{\Delta}(a \circ b)$$

が成立する. $\mathbf{X} = [a, b]$, $\mathbf{Y} = [c, d]$ ($a, b, c, d \in \mathbb{F}$) に対して, 機械区間演算は次のように実現できる.

$$\mathbf{X} + \mathbf{Y} = [\text{fl}_{\nabla}(a + c), \text{fl}_{\Delta}(b + d)]$$

$$\mathbf{X} - \mathbf{Y} = [\text{fl}_{\nabla}(a - d), \text{fl}_{\Delta}(b - c)]$$

$$\mathbf{X} \times \mathbf{Y} = [\text{fl}_{\nabla}(\min\{ac, ad, bc, bd\}), \text{fl}_{\Delta}(\max\{ac, ad, bc, bd\})]$$

$$\mathbf{X} \div \mathbf{Y} = [\text{fl}_{\nabla}(\min\{a/c, a/d, b/c, b/d\}), \text{fl}_{\Delta}(\max\{a/c, a/d, b/c, b/d\})]$$

$\mathbf{X} \times \mathbf{Y}$	$c > 0$	$0 \in \mathbf{Y}$	$d < 0$
$a > 0$	$[\text{fl}_{\nabla}(ac), \text{fl}_{\Delta}(bd)]$	$[\text{fl}_{\nabla}(bc), \text{fl}_{\Delta}(bd)]$	$[\text{fl}_{\nabla}(bc), \text{fl}_{\Delta}(ad)]$
$0 \in \mathbf{X}$	$[\text{fl}_{\nabla}(ad), \text{fl}_{\Delta}(bd)]$	B	$[\text{fl}_{\nabla}(bc), \text{fl}_{\Delta}(ac)]$
$b < 0$	$[\text{fl}_{\nabla}(ad), \text{fl}_{\Delta}(bc)]$	$[\text{fl}_{\nabla}(ad), \text{fl}_{\Delta}(ad)]$	$[\text{fl}_{\nabla}(bd), \text{fl}_{\Delta}(ac)]$

ただし $B = [\min\{\text{fl}_{\nabla}(ad), \text{fl}_{\nabla}(bc)\}, \max\{\text{fl}_{\Delta}(ad), \text{fl}_{\Delta}(bc)\}]$.

2.5. ベクトル・行列の区間演算

上で述べた丸めの向きを制御することにより、ベクトル $x, y \in \mathbb{F}^n$ の内積 $x^T y$ 、行列 $A, B \in \mathbb{F}^{n \times n}$ の積、あるいは、ベクトル行列積 Ax の結果を区間で厳密に包含することができる。

$$\begin{aligned} \text{fl}_{\nabla}(x^T y) &\leq x^T y \leq \text{fl}_{\Delta}(x^T y) \\ \text{fl}_{\nabla}(Ax) &\leq Ax \leq \text{fl}_{\Delta}(Ax) \\ \text{fl}_{\nabla}(AB) &\leq AB \leq \text{fl}_{\Delta}(AB) \end{aligned}$$

このようにすると丸め方向の制御で区間演算が容易にできる。しかし、行列ベクトル積、行列積を高速に実装することは職人芸のレベルの難しさである（例えば、キャッシュサイズをみて最適なブロック分割などを行う）。そのため通常は数値計算ライブラリを利用するのが主流である。

演習問題

3つの浮動小数点数 $a, b, c \in \mathbb{F}$ に対して、次が正しいかどうかを調べよ。正しくない場合は理由を述べ、できたらその反例をあげよ。

1. $\text{fl}_{\Delta}(a + b) \geq a + b$
2. $\text{fl}_{\Delta}(a - b) \geq a - b$
3. $\text{fl}_{\Delta}(ab) \geq ab$
4. $\text{fl}_{\Delta}(a/b) \geq a/b$
5. $\text{fl}_{\Delta}((a + b) + c) \geq (a + b) + c$
6. $\text{fl}_{\Delta}((ab)c) \geq (ab)c$
7. $\text{fl}_{\Delta}((a - b) - c) \geq (a - b) - c$
8. $\text{fl}_{\Delta}(a - (b + c)) \geq a - (b + c)$
9. $\text{fl}_{\Delta}(-a + b) \geq (-a + b)$
10. $\text{fl}_{\Delta}(-((-a) + (-b))) \geq -((-a) + (-b))$
11. $\text{fl}_{\Delta}(ab + cd) \geq ab + cd$
12. $\text{fl}_{\Delta}((a + b)(c + d)) \geq (a + b)(c + d)$
13. $\text{fl}_{\Delta}(a/(b + c)) \geq a/(b + c)$
14. $\text{fl}_{\Delta}(a - bc) \geq a - bc$
15. $\text{fl}_{\Delta}(a + (-b)c) \geq a + (-b)c$

16. $f1_{\Delta}(\text{sqrt}(a)) \geq \text{sqrt}(a)$