

浮点数的表示与算术运算算法

姓名：应承峻

学号：3170103456

实验描述

浮点数的表示与算术运算算法分析，要求理论推导与程序模拟。算术运算包括字符串到浮点数、浮点数到字符串的转换，加、减、乘、除四则运算等。要有说明文档，包括算法证明、程序框图、使用方法、特殊处理(溢出、数位扩展)、实例分析等等。字符串转换可能稍难。

```
1 typedef unsigned int dwrd;
2 char* ftoa(dwrd);
3 dwrd atof(char*);
4 dwrd fadd(dwrd, dwrd);
5 dwrd fsub(dwrd, dwrd);
6 dwrd fmul(dwrd, dwrd);
7 dwrd fdiv(dwrd, dwrd);
```

- 主程序只供验证，你所须写的就是6个子程序。atof 和 ftoa 可能会比较麻烦。
- 对结果应该进行分析，讨论你的浮点运算适用范围，可能的问题等等。
- 验证要有大数、小数，比如：12345678901234567890123456789.123456789。浮点的范围是十的正负38次方。如果不能，分析原因。
- 可以增加十六进制的对比。比如显示 z1.f 与 (u.f+v.f) 的十六进制。这样可以分析四舍五入情况。

算法分析

加减法

浮点数加法的原则为：“小阶向大阶看齐”，根据浮点数的格式取出 ex, ey 分别为浮点数 x 和 y 的阶码，而 mx, my 为其尾数。计算时，将小阶对应的尾数右移后向大阶看齐，然后将尾数相减得到新的尾数。相减后进行重新规格化，但此时要注意如果两尾数相减后为0则直接返回0即可否则会导致死循环。

浮点数的减法可以通过加法来实现： $x - y = x + (-y)$ ，因此只需改变 y 的符号位后再与 x 做加法即可。

```
1 dwrd fadd(dwrd x , dwrd y) {
2     if (x == 0) return y;
3     else if (y == 0) return x;
4     if ((x << 1) < (y << 1)) return fadd(y , x); // |x| >= |y|
5     dwrd ex = (x >> 23) & 0xFF;
6     dwrd ey = (y >> 23) & 0xFF;
7     dwrd mx = (x & 0x7FFFFFFF) | 0x800000;
8     dwrd my = (y & 0x7FFFFFFF) | 0x800000;
```

```

9   my >>= (ex - ey);
10  if ((x^y) & 0x80000000) { //case -
11      mx -= my;
12      if (mx == 0) return 0; //zero
13      while ((mx & 0x00800000) == 0) {
14          mx <<= 1;
15          ex--;
16      }
17  } else { //case +
18      mx += my;
19      if (mx == 0) return 0; //zero
20      while (mx & 0xFF000000) {
21          mx >>= 1;
22          ex++;
23      }
24  }
25  return (x & 0x80000000) | (ex << 23) | (mx & 0x7FFFFFFF);
26 }
27
28 dwrdrd fsub(dwrdrd x , dwrd y) {
29     return fadd(x , y ^ 0x80000000);
30 }

```

乘法

乘法的阶码运算法则，设 $E(x)$ 为阶码， x 为阶：

$$E(x) = x + M$$

$$E(xy) = x + y + M = E(x) + E(y) - M$$

根据真值表，浮点数乘法的符号位即为若两数最高位异或值：

| x | y | $(x \oplus y) \wedge 0x80000000$ |
|-----|-----|----------------------------------|
| 1 | 0 | 1 |
| 1 | 1 | 0 |
| 0 | 1 | 1 |
| 0 | 0 | 0 |

乘法的尾数即为两浮点数尾数值相乘，但由于一个浮点数的尾数（此处加上最前面的1）有24位，两个浮点数相乘时达48位，导致发生溢出。所以，程序中只取两个尾数相乘后的前24位作为新的尾数而舍弃了后24位。

浮点数乘法后可能会产生类似加法一样的进位问题，所以在得到结果后也要进行规格化处理。

```

1  dwrd fmul(dwrdrd x, dwrd y) {
2      if (x == 0 || y == 0) return 0;
3      dwrd ex = (x >> 23) & 0xFF;
4      dwrd ey = (y >> 23) & 0xFF;
5      dwrd mx = (x & 0x7FFFFFFF) | 0x800000; //24-bit
6      dwrd my = (y & 0x7FFFFFFF) | 0x800000;

```

```

7   dwrd mr = 0;
8   ex = ex + ey - 127;
9   while (my != 0) {    //乘法算法
10      mr >>= 1;
11      if (my & 1) mr += mx;
12      my >>= 1;
13  }
14  while (mr & 0xFF000000) { //规格化
15      mr >>= 1;
16      ex++;
17  }
18  return ((x^y) & 0x80000000) | (ex << 23) | (mr & 0x7FFFFFFF);
19 }

```

除法

除法的阶码运算法则，设 $E(x)$ 为阶码， x 为阶：

$$E(x/y) = x - y + M = E(x) - E(y) + M$$

尾数除法采用的算法为恢复除数的除法，为了便于对算法的理解我们以 $4 \div 10$ 为例进行演示

4.0在浮点数中表示为：0 10000001 000000000000000000000000

10.0在浮点数中表示为：0 10000010 010000000000000000000000

首先通过阶码的运算法则得到阶码为 $10000001 - 10000010 + 01111111 = 01111110$ (126)

然后对尾数做除法 $100000000000000000000000 \div 101000000000000000000000$

为了处理最高位商0导致的移位失败，我们在开始时通过不断的移位处理使得 $mx > my$ ，这样保证了最高位能够商1，除此之外，该操作也能够保证除法后不需要再进行借位规格化处理。此外，该操作也能排除左移导致的进位规格化处理的问题，证明如下：设最后一次移位前 $x/y = \alpha$ ，由移位条件得 $\alpha \in [1/2, 1)$ ，当再进行移位后，可以得到 $\alpha_1 \in [1, 2)$ ，因此商必介于1与2之间，不会产生进位。

每次除法时，如果 $mx < my$ ，则说明不够减，此时商0，否则将 mx 减去 my 并且商1，最后每次减法后将被除数左移，直到商的第24位为1为止。

```

1   dwrd fdiv(dwrd x , dwrd y) {
2       if (x == 0) return 0;
3       if (y == 0) {
4           cout << "除数不能为0! " << endl;
5           return (x & 0x80000000) | 0x7F800000; //INF
6       }
7       dwrd ex = (x >> 23) & 0xFF;
8       dwrd ey = (y >> 23) & 0xFF;
9       dwrd mx = (x & 0x7FFFFFFF) | 0x800000;
10      dwrd my = (y & 0x7FFFFFFF) | 0x800000;
11      dwrd mr = 0;
12      ex = ex - ey + 127;
13      while (mx < my) { //handle borrow before calculate
14          mx <<= 1;
15          ex--;
16      }
17      while ((mr & 0xFF800000)==0) {

```

```

18     if (mx < my) mr <= 1; //quotient 0
19     else {
20         mx -= my;
21         mr = (mr << 1) | 1; //quotient 1
22     }
23     mx <= 1;
24 }
25 return ((x^y) & 0x80000000) | (ex << 23) | (mr & 0x7FFFFF);
26 }

```

浮点数转换成字符串

浮点数转换成字符串较为繁琐，首先需要考虑特殊情况，如果是0则直接输出0，如果阶码是255则直接输出INF表示无穷大。其次是考虑符号位，根据最高位判断出符号位后直接加到字符串中。

接下来需要分情况讨论： x 代表一个大于1的数还是小于1的数。

如果是一个大于1的数，则需要得到其整数部分和小数部分。

以12(1100.00000000000000000000)为例

整数部分较为容易，用高精度算法直接乘幂即可。而小数部分的处理比较麻烦。由于 $x \div 2 = x \times 5/10$ ，因此我们可以通过“补0乘5”的方式得到每次运算的权重，如 $0.5 \rightarrow 0.25 \rightarrow 0.125 \rightarrow 0.0625$ 的权重变化可以看做是：

$5 \rightarrow 05 * 5 = 25 \rightarrow 025 * 5 = 125 \rightarrow 0125 * 5 = 0625 \rightarrow \dots$ 这一变换得到的，因此我们每次构造一个以5位幂的权重进行计算即可。

如果 x 是一个小于1的数，则不需要处理整数部分，小数部分的处理方法同上。

```

1  string ftoa(dword x) {
2      if (x == 0) return "0"; //special case zero
3      string buffer;
4      if (x & 0x80000000) buffer += '-'; //x<0
5      dword ex = (x >> 23) & 0xFF; //exp
6      dword mx = (x & 0x7FFFFFFF) | 0x800000; //mx
7      dword integer = 0, decimal = 0;
8      if (ex == 255) buffer += "INF"; //infinity
9      else if (ex >= 127) {
10         ex -= 127;
11         string result = "0", power = "1";
12         for (int k = ex + 1; k > 0; k--) {
13             if ((mx >> (24 - k)) & 1) result = add(result, power);
14             power = multi(power, 2);
15         }
16         buffer = buffer + result + ".";
17         result = "0"; power = "5";
18         for (int k = ex + 2; k <= 24; k++) {
19             if ((mx >> (24 - k)) & 1) {
20                 string temp = power;
21                 reverse(temp.begin(), temp.end());
22                 result = add(result, temp);
23             }
24             power = "0" + power;
25             power = multi(power, 5);

```

```

26     }
27     reverse(result.begin() , result.end());
28     buffer += result;
29 } else { //<1
30     buffer += "0.";
31     int len = 24;
32     string result="0" , power = "5";
33     for (dwrdr k = 1; k < 127 - ex; k++) {
34         power = "0" + power;
35         power = multi(power , 5);
36     }
37     while ((mx & 1) == 0) {
38         mx >>= 1;
39         len--;
40     }
41     while (len > 0) {
42         if ((mx >> (len - 1)) & 1) {
43             string temp = power;
44             reverse(temp.begin() , temp.end());
45             result = add(result , temp);
46         }
47         power = "0" + power;
48         power = multi(power , 5);
49         len--;
50     }
51     reverse(result.begin() , result.end());
52     buffer += result;
53 }
54 return buffer;
55 }

```

字符串转换成浮点数

在将字符串转换成浮点数前，首先通过正则表达式检验给定字符串是否是一个合法的浮点数。

为了能够分别计算整数部分和小数部分，我们将整个字符串拆分成三部分：符号位、整数子串和小数子串。拆分后，需要考虑一个特殊情况：即浮点数为0的情况，此时应直接返回。

整数部分的转化思路为除二取余，以此得到整数部分的二进制串。小数部分的转换比较麻烦，我们以0.1234为例子，通过字符串切割得到小数部分1234，乘二取整的结果取决于是否前面多出了一位，比如5001乘二时得到了10002，此时需要取证，而 $1234 \times 2 = 2468 < 10000$ 不需要取证，因此设小数部分的字符串长度为 `dotsize`，则乘二取整得到1的条件为 `multi(decimal,2).size()>=dotsize`。

将整数和小数的字符串拼接后截取前23位作为尾数，再与符号位和阶码拼接即可。

```

1  dwrdr atof(string s) {
2      bool borrow = false;
3      bool check = regex_match(s , regex("^[-]?[0-9]*\\.?[0-9]+$"));
4      if (!check) {
5          cout << "不是一个合法的浮点数" << endl;
6          return 0;
7      }
8      string bin , integer = "0" , decimal = "0" , res;

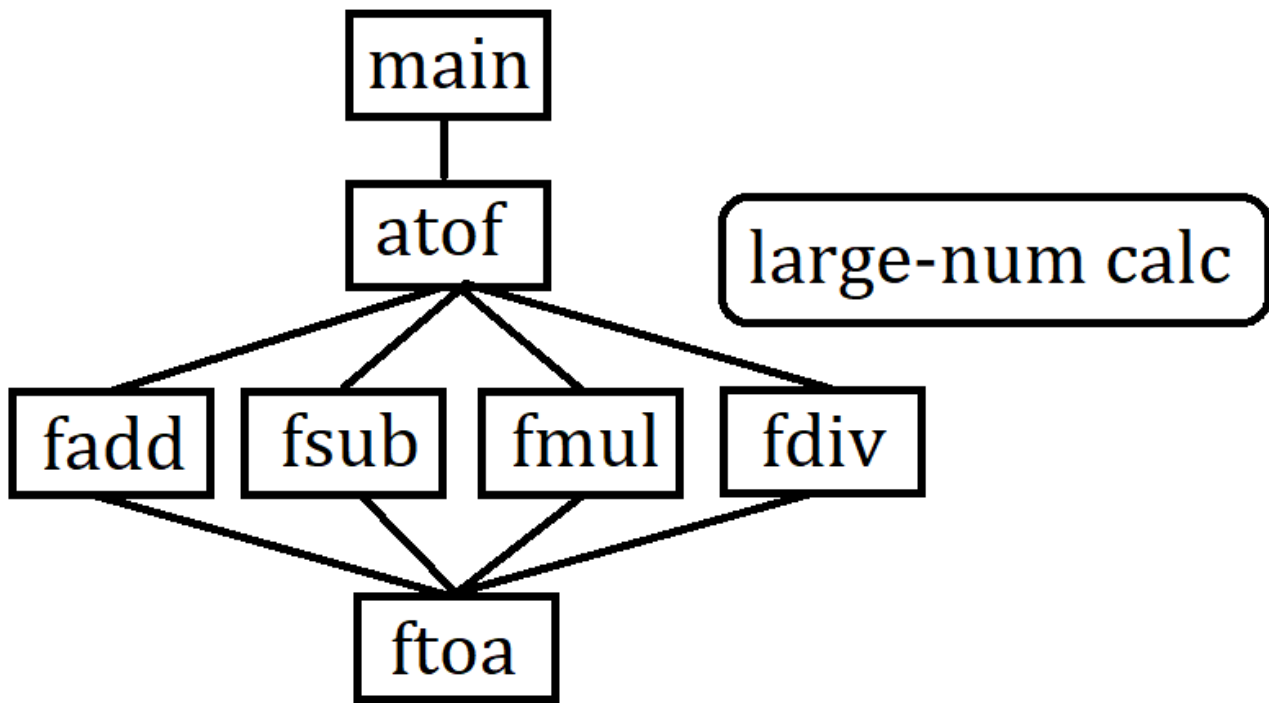
```

```

9      dwrd pos , i = 0;
10     dwrd ex = 0 , x = 0;
11     if (s.at(0) == '-') { //calc symbol
12         res += '1';
13         i++;
14     } else res += '0';
15     res += "00000000";
16     pos = i;
17     while (pos != s.size() && s.at(pos) != '.' ) pos++; //seperate
18     if (pos < s.size()) decimal = s.substr(pos + 1);
19     integer = s.substr(i , pos - i);
20     if (!integer.compare("0") && !decimal.compare("0")) return 0; //special case for
input zero
21     if (integer.compare("0") == 0) bin = "0";
22     else {
23         while (integer.compare("0") > 0) { //convert to binary
24             bin += (mod(integer , 2)) ? "1" : "0";
25             integer = divide(integer , 2);
26         }
27         reverse(bin.begin() , bin.end());
28     }
29     ex = bin.size() - 1 + 127; //evaluate ex
30     dwrd dotsize = decimal.size();
31     while (bin.size() < 46) { //construct 23-bit format mx
32         decimal = multi(decimal , 2);
33         if (decimal.size() > dotsize) {
34             decimal.erase(decimal.begin());
35             bin += "1";
36         } else bin += "0";
37     }
38     while (bin.at(0) == '0') { //case for type 0.***
39         bin.erase(bin.begin());
40         ex--;
41     }
42     bin.erase(bin.begin());
43     res += bin.substr(0 , 23);
44     for (i = 0; i < 32; i++) x = x * 2 + (res.at(i) - '0');
45     return x + (ex << 23);
46 }

```

程序框图



特殊处理

针对 `unsigned int` 在字符串与浮点数相互转换的溢出情况，我们采用了高精度算法进行计算使其能够达到数位扩展的要求，其算法如下。

```
1  string divide(string s , dwr d b) { //large-number divide algorithm
2      dwr d r = 0;
3      dwr d i = 0;
4      char c[100];
5      char *p;
6      string res;
7      for (i = 0; i < s.size(); i++) {
8          r = r * 10 + s[i] - '0';
9          if (r < b) c[i] = '0';
10         else {
11             c[i] = r / b + '0';
12             r = r % b;
13         }
14     }
15     c[i] = '\0';
16     p = c;
17     while ((*p == '0') && p) p++;
18     return p;
19 }
20
21 dwr d mod(string s , dwr d b) { //large-number mod algorithm
22     dwr d r = 0;
23     dwr d i = 0;
24     char c[100];
25     char *p;
```

```

26     string res;
27     for (i = 0; i < s.size(); i++) {
28         r = r * 10 + s[i] - '0';
29         if (r < b) c[i] = '0';
30         else {
31             c[i] = r / b + '0';
32             r = r % b;
33         }
34     }
35     c[i] = '\0';
36     p = c;
37     while ((*p == '0') && p) p++;
38     return r;
39 }
40
41
42 string multi(string s , dwrd b) { //large-number multiple algorithm
43     char c[100];
44     dwrd len = 0;
45     dwrd carry = 0;
46     for (int i = s.size() - 1; i >= 0; i--) {
47         dwrd temp = (s[i] - '0') * b + carry;
48         c[len++] = temp % 10 + '0';
49         carry = temp / 10;
50     }
51     while (carry != 0) {
52         c[len++] = carry % 10 + '0';
53         carry /= 10;
54     }
55     c[len] = '\0';
56     string res = c;
57     reverse(res.begin() , res.end());
58     return res;
59 }
60
61
62 string add(string x , string y) { //large-number add algorithm
63     char c[100];
64     string res;
65     dwrd len = 0;
66     dwrd carry = 0;
67     reverse(x.begin() , x.end());
68     reverse(y.begin() , y.end());
69     while (x.size() < y.size()) x += '0';
70     while (y.size() < x.size()) y += '0';
71     for (int i = 0; i < x.size(); i++) {
72         int temp = x[i] - '0' + y[i] - '0' + carry;
73         c[len++] = temp % 10 + '0';
74         carry = temp / 10;
75     }
76     if (carry != 0) {
77         c[len++] = carry + '0';
78     }

```



```

79     c[len] = '\0';
80     res = c;
81     reverse(res.begin() , res.end());
82     return res;
83 }

```

使用方法与实例分析

使用方法：根据提示输入两个浮点数即可。

测试样例1：能够转换成二进制的浮点数，带符号

```

Float 1: 2.5
Float 2: -4
Float 1 is : 2.5
Float 2 is : -4
Sum(x, y) = -1.5
Sub(x, y) = 6.5
Mul(x, y) = -10
Div(x, y) = -0.625
请按任意键继续. . .

```

测试样例2：不能够准确转换成二进制的浮点数，带符号

```

Float 1: 2.4
Float 2: -3.1
Float 1 is : 2.4
Float 2 is : -3.1
Sum(x, y) = -0.7
Sub(x, y) = 5.5
Mul(x, y) = -7.44
Div(x, y) = -0.774194
请按任意键继续. . .

```

测试样例3：大数运算（分两行，前面一行为标准输出，后一行为atof函数输出）

```

Float 1: 123456789
Float 2: 876543210
Float 1 is : 1.23457e+08
Float 2 is : 8.76543e+08
Sum(x, y) = 1e+09
Sum(x, y) = 999999936.0
Sub(x, y) = -7.53086e+08
Sum(x, y) = -753086400.0
Mul(x, y) = 1.08215e+17
Mul(x, y) = 108215197151294758.0
Div(x, y) = 0.140845
Div(x, y) = 0.1308339493396020512975
请按任意键继续. . .

```

测试样例4：大数运算（分两行，前面一行为标准输出，后一行为atof函数输出）

```
Float 1: 12345678901234567890123456789.123456789
Float 2: -2.33333333333333
Float 1 is : 1.23457e+28
Float 2 is : -2.33333
Sum(x, y) = 1.23457e+28
Sum(x, y) = 12345677976674082762172668672.0
Sub(x, y) = 1.23457e+28
Sum(x, y) = 12345677976674082762172668672.0
Mul(x, y) = -2.88066e+28
Mul(x, y) = -28806579584389617793824472704.0
Div(x, y) = -5.291e+27
Div(x, y) = -5291004931474008397523230176.0
请按任意键继续. . .
```

分析：本程序浮点数适用范围：由于采用了高精度算法，因此在10的正负38次方内均可以使用。但是由于 `atof` 函数在转换输出时存在一定的误差，因此针对不能完全转换成二进制浮点数的数的输出会存在精度丢失的问题。