

计算物理学第一次作业报告

陈跃元

物理学院 1600011328

1 双精度浮点数和其 IEEE754 编码的互相转化

1.1 问题描述

实现任意十进制数于其规范化的二进制双精度浮点数编码间的相互转换. 具体而言, 输入一个双精度浮点数, 需要给出它按照 IEEE754 标准的编码; 输入一个 64 位, 由 ‘0’ 和 ‘1’ 组成的字符串, 将它按照 IEEE754 标准转化为双精度浮点数.

1.2 算法

1.2.1 双精度浮点数转编码

有两种可行的算法:

1. 按定义转换

- 确定符号位 (若 $x == 0$ 直接返回全零), exponent 初始值设为 1023.
- (确定 exponent) 判断 $x \geq 1$?
 - Yes
不断重复 $x /= 2$, exponent ++, 直到 $x < 2$.
 - No
不断重复 $x *= 2$, exponent --, 直到 $x \geq 1$.
- $x --$
- (小数的二进制转换) 从 $i = 12$ 开始, 不断重复 $x *= 2$, $s[i] = x \% 1$, $i ++$, 直至将 s 的小数部分 52 位全部填满.
- 返回 s .

这一算法比起将 x 分解为整数部分和小数部分分别转换的效率要更高 (同为 $O(\lg n)$ 级别, 但常数比分开转换要小), 且在大数情形下, 不会出现分解出的整数部分过大造成的 overflow.

2. 利用强制指针类型转换

这一算法某种程度上是一种投机取巧的办法. 在 C++ 中, double 类型和 long long 类型都是由 8 个字节来存储的, 且 double 类型在内存中的存储方式就是按照 IEEE 754 标准执行的, 故可以通过将指向 x 的 double 类指针强制转化为 long long 类指针, 随后将这一 long long 按照整数的二进制转换来得到编码.

1.2.2 编码转双精度浮点数

这时的算法相当于将上面两个算法倒过来.

1. 按定义转换

- (确定小数) 将 s 的后 52 位按照二进制小数的规则转化为 x .
- $x++$
- (确定 exponent) 将 exponent 部分计算出来. 判断:
 - exponent == 0
返回 $x = 0$.
 - exponent >= 1023
不断重复 $x* = 2$, exponent --, 直到 exponent = 1023.
 - $0 < \text{exponent} < 1023$
不断重复 $x/ = 2$, exponent ++, 直到 exponent = 1023.
- 若符号位为 1, $x* = -1$.
- 返回 x .

2. 利用强制指针类型转换

将 s 按照整数的二进制转化规则转化为 long long 类型的 x , 并将指向 x 的 long long 指针强制类型转换为 double 指针, 返回指针所指的双精度浮点数.

源代码见附录A.1.1, 或 src/Problem_1/main.cpp.

1.3 结果

```
Example from Chapter 1 Lecture note Page 50
Convert -222.111 by method 1:
1100000001101011110000111000110101001111110111110011101101100100
Convert -222.111 by method 2:
1100000001101011110000111000110101001111110111110011101101100100
Convert back by method 1:
-222.111
Convert back by method 2:
-222.111
Try another example. Input the double:
1234.567890
Convert 1234.56789 by method 1:
0100000010010011010010100100010110000100111101001100011011100111
Convert 1234.56789 by method 2:
0100000010010011010010100100010110000100111101001100011011100111
Convert back by method 1:
1234.56789
Convert back by method 2:
1234.56789
Try another example. Input the string:
10111011010111110111101000011110010101010101011110001111111110
Convert back by method 1:
-1.04148740441655e-22
Convert back by method 2:
-1.04148740441655e-22
```

图 1: 第一题的结果

运行如图1所示. 两种方法得到的结果相同.

2 带权的高斯积分

2.1 问题描述

构造下列三个高斯积分的节点和权重.

$$\int_0^1 \sqrt{x} f(x) dx \approx A_0 f(x_0) + A_1 f(x_1) \quad (1)$$

$$\int_0^1 \sqrt{x} f(x) dx \approx A_0 f(x_0) + A_1 f(x_1) + A_2 f(x_2) \quad (2)$$

$$\int_{-1}^1 (1+x^2) f(x) dx \approx A_0 f(x_0) + A_1 f(x_1) \quad (3)$$

2.2 推导和结果

对 (1), 对于节点 x_i , 应满足对所有 $0 \leq k \leq n$,

$$\int_0^1 \sqrt{x} \left[\prod_{i=0}^n (x - x_i) \right] P_k(x) dx = 0 \quad (4)$$

即

$$\begin{cases} \frac{1}{7} - \frac{x_0}{5} - \frac{x_1}{5} + \frac{x_0 x_1}{3} = 0 \\ \frac{1}{9} - \frac{x_0}{7} - \frac{x_1}{7} + \frac{x_0 x_1}{5} = 0 \end{cases} \quad (5)$$

解得 $x_0 = \frac{1}{63}(35 - 2\sqrt{70}) \approx 0.289949$, $x_1 = \frac{1}{63}(35 + 2\sqrt{70}) \approx 0.821162$. 带入原式, 为使所有阶数 ≤ 1 的多项式积分为准确值, 应有

$$\begin{cases} A_0 + A_1 = \frac{2}{3} \\ A_0 x_0 + A_1 x_1 = \frac{2}{5} \end{cases} \quad (6)$$

解得 $A_0 \approx 0.277556$, $A_1 \approx 0.389111$.

对 (2), 同样有 (4) 的要求, 即

$$\begin{cases} -\frac{2}{3}abc + \frac{2}{5}(ab + bc + ac) - \frac{2}{7}(a + b + c) + \frac{2}{9} = 0 \\ -\frac{2}{5}abc + \frac{2}{7}(ab + bc + ac) - \frac{2}{9}(a + b + c) + \frac{2}{11} = 0 \\ -\frac{2}{21}abc + \frac{2}{15}(ab + bc + ca) - \frac{10}{77}(a + b + c) + \frac{14}{117} = 0 \end{cases} \quad (7)$$

解得 $x_0 = \frac{1}{63}(35 - 2\sqrt{70}) \approx 0.289949$, $x_1 = \frac{1}{63}(35 + 2\sqrt{70}) \approx 0.821162$. 带入原式, 为使所有阶数 ≤ 1 的多项式积分为准确值, 应有

$$\begin{cases} A_0 + A_1 = \frac{2}{3} \\ A_0 x_0 + A_1 x_1 = \frac{2}{5} \end{cases} \quad (8)$$

解得 $A_0 \approx 0.277556$, $A_1 \approx 0.389111$.

3 衰变

3.1 问题描述

两类核 A 和 B 发生衰变, 初始时刻 $N_A = N_B = 1$, 满足

$$\frac{dN_A}{dt} = -\frac{N_A}{\tau_A} \quad (9)$$

$$\frac{dN_B}{dt} = \frac{N_A}{\tau_A} - \frac{N_B}{\tau_B} \quad (10)$$

编写程序求解耦合微分方程组, 并和解析解比较.

3.2 算法

解析解:

$$N_A = \exp(-\frac{t}{\tau_A}), \quad N_B = \frac{\tau_A \exp(-\frac{t}{\tau_B}) + \tau_B (\exp(-\frac{t}{\tau_A}) - 2 \exp(-\frac{t}{\tau_B}))}{\tau_A - \tau_B} \quad (11)$$

笔者考虑了下面四种不同的算法 (伪代码中将 N_A 和 N_B 分别记作 a, b):

1. Naive iteration

```
1      a[n] = a[n-1] - delta * a[n-1] / tau_a
2      b[n] = b[n-1] + delta * (a[n-1] / tau_a - b[n-1] / tau_b)
```

其中 n 是迭代步数, delta 是时间的步长, tau_a 和 tau_b 的意义不言自明, 下同.

2. Naive iteration, 但对 N_B 迭代时, N_A 取前后两次的 N_A 的平均值.

```
1      a[n] = a[n-1] - delta * a[n-1] / tau_a
2      b[n] = b[n-1] + delta * ((a[n] + a[n-1]) / (2 * tau_a) - b[n-1] / tau_b)
```

3. 对 N_A 用 RK-4, 对 N_B 用 Naive iteration, 但 N_A 的值取两次迭代的平均值.

```
1      k[0] = - delta * a[n-1] / tau_a
2      k[1] = - delta * (a[n-1] + k[0] / 2) / tau_a
3      k[2] = - delta * (a[n-1] + k[1] / 2) / tau_a
4      k[3] = - delta * (a[n-1] + k[2]) / tau_a
5      a[n] = a[n-1] + (k[0] + 2 * k[1] + 2 * k[2] + k[3]) / 6
6      b[n] = b[n-1] + delta * ((a[n-1] + a[n]) / (2 * tau_a) - b[n-1] / tau_b)
```

4. 对 N_A 用 RK-4, 然后再对 N_B 用 RK-4.

```
1      k[0] = -delta * a[n-1] / tau_a
2      k[1] = -delta * (a[n-1] + k[0] / 2) / tau_a
3      k[2] = -delta * (a[n-1] + k[1] / 2) / tau_a
4      k[3] = -delta * (a[n-1] + k[2]) / tau_a
5      a[n] = a[n-1] + (k[0] + 2 * k[1] + 2 * k[2] + k[3]) / 6
6
7      k[0] = delta * (a[n-1] / tau_a - b[n-1] / tau_b);
8      k[1] = delta * ((a[n-1] + a[n]) / (2 * tau_a) - (b[n-1] + k[0] / 2) / tau_b)
9      // Explicit time dependent in N_A, linear intrapolate
10     k[2] = delta * ((a[n-1] + a[n]) / (2 * tau_a) - (b[n-1] + k[1] / 2) / tau_b)
11     // Explicit time dependent in N_A, linear intrapolate
12     k[3] = delta * (a[n] / tau_a - (b[n-1] + k[2]) / tau_b)
13     b[n] = b[n-1] + (k[0] + 2 * k[1] + 2 * k[2] + k[3]) / 6
```

需要注意一点, 在得到 $a[n]$ 以后, 代入计算 N_B 时, RK-4 要求 $k[1]$ 和 $k[2]$ 中取 $t[n] + \text{delta} / 2$ 时的 N_A , 但这一时刻并不在迭代的时刻中, 考虑到 delta 较小, 期间 N_A 的变化可以近似看做线性的, 故通过线性插值得到 $t[n] + \text{delta} / 2$ 时的 N_A .

源代码见附录A.1.2, 或 `src/Problem_3/main.cpp`.

3.3 结果

3.3.1 几种方法的比较

笔者首先比较了不同算法在 $\tau_A = 1\text{s}, \tau_B = 10\text{s}, \Delta t = 0.2\text{s}$ 的情况下的误差, 以找到其中误差最小的方法. 图2、3、4分别给出了各个方法的计算结果、相对误差和绝对误差. 表1中列出了各种方法的最大相对误差. 很显然, 第四种算法的误差相比于剩下三种方法的误差要小得多, 故之后的分析都基于第四种方法.

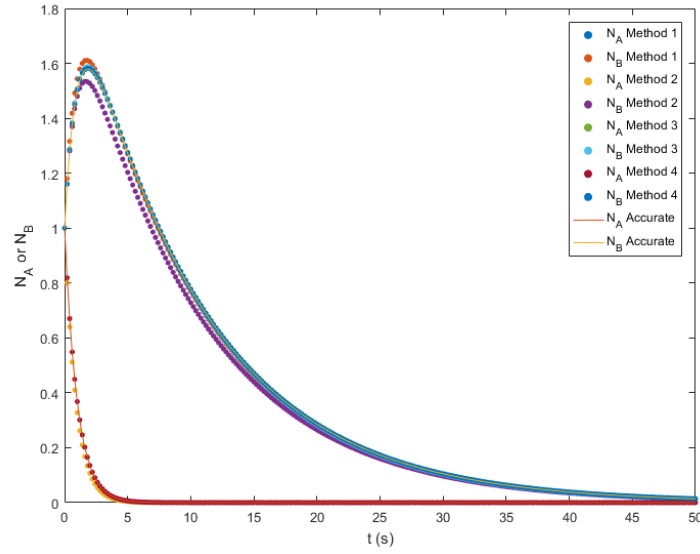


图 2: 不同方法在 $\tau_A = 1\text{s}, \tau_B = 10\text{s}, \Delta t = 0.2\text{s}$ 的情况下的结果

表 1: 不同方法在 $\tau_A = 1\text{s}, \tau_B = 10\text{s}, \Delta t = 0.2\text{s}$ 的情况下的最大相对误差 (绝对值)

	Method 1	Method 2	Method 3	Method 4
N_A	99.69%	99.69%	0.10%	0.10%
N_B	4.95%	9.94%	4.18%	0.22%

3.4 不同的 τ_B 下的行为

如图5所示, N_A 的行为与 τ_A 和 τ_B 的大小关系无关, 而当 $\tau_B > \tau_A$ 时, N_B 短期行为将会出现增加, $\tau_B \leq \tau_A$ 时, N_B 将一直减小. N_B 的长期行为成指数衰减.

3.5 步长对误差的影响

表2给出了 $t = 50\text{s}$ 内不同步长 Δt 下的最大相对误差 (绝对值). 由于算法本身的误差已经很小, 步长变化造成的影响很小, 但依然可以看出, 步长越小, 误差越小.

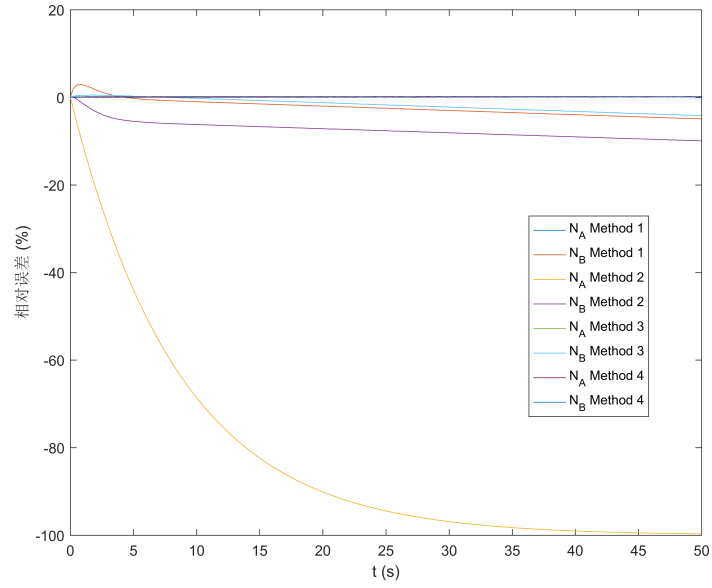


图 3: 不同方法在 $\tau_A = 1\text{ s}$, $\tau_B = 10\text{ s}$, $\Delta t = 0.2\text{ s}$ 的情况下的相对误差

表 2: $t = 50\text{ s}$ 内不同步长 Δt 下的最大相对误差 (绝对值)

Δt	0.2 s	0.1 s	0.05 s
N_A	0.103%	0.048%	0.048%
N_B	0.219%	0.089%	0.057%

4 双精度浮点数和其 IEEE754 编码的互相转化

4.1 问题描述

4.2 算法

4.3 结果

4.4 源代码

5 双精度浮点数和其 IEEE754 编码的互相转化

5.1 问题描述

5.2 算法

5.3 结果

5.4 源代码

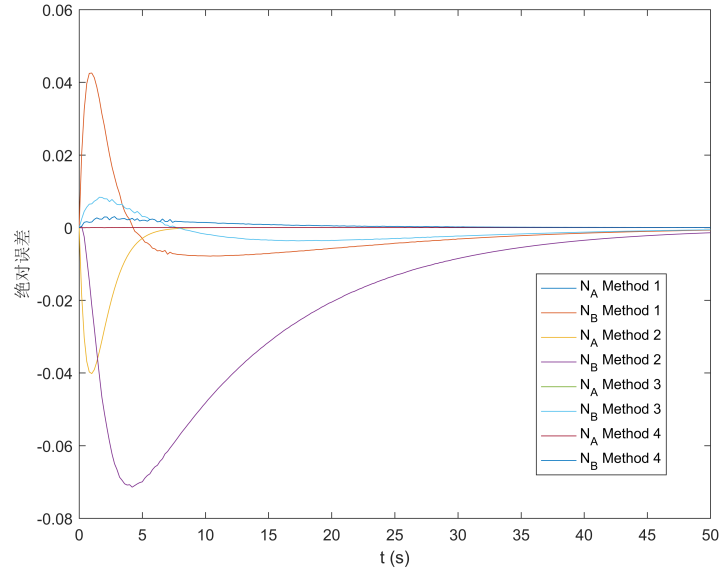


图 4: 不同方法在 $\tau_A = 1\text{ s}$, $\tau_B = 10\text{ s}$, $\Delta t = 0.2\text{ s}$ 的情况下的绝对误差

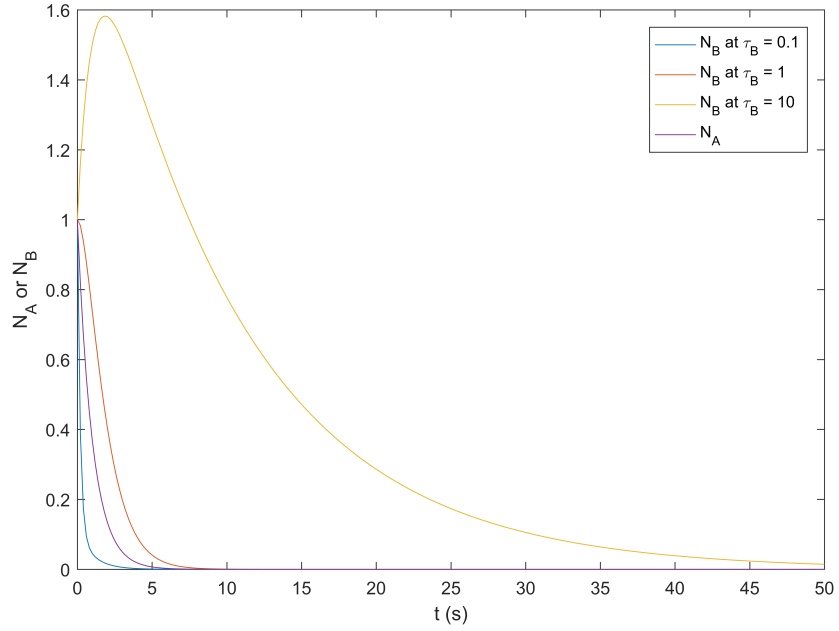


图 5: 不同 τ_B 下的结果

A 附录

A.1 源代码

A.1.1 第一题

```
1 #include <iostream>
2 #include <string>
3 #include <iomanip>
4
5 using namespace std;
6
7 /*
8 =====
9 Program Description
10 -----
11 This program contains the following functions:
12
13 1. char * double_2_ieee_method_1:
14    This method converts a double to a c string of its IEEE 754 format, by using
15    the definition of double-precision floating point format of IEEE 754 2008
16    standard.
17
18 2. double ieee_2_double_method_1:
19    This method converts a c string of IEEE 754 format to a double, by using the
20    definition of double-precision floating point format of IEEE 754 2008
21    standard.
22
23 3. char * double_2_ieee_method_2:
24    This method converts a double to a c string of its IEEE 754 format, by
25    explicitly converting the double to a long long, reading the converted
26    number bit by bit, and writing these bits into the string.
27
28 4. double ieee_2_double_method_2:
29    This method converts a c string of IEEE 754 format to a double, by
30    converting the string to a long long and then explicitly convert the long
31    long to a double.
32
33 The other methods are for presentation of the results.
34
35 Detailed explanation of the algorithm and the results can be found in the
36 report.
37 =====
38 */
39
40 // Convert by definition (IEEE 754)
41 // ** Subnormal numbers are NOT considered **
42 char * double_2_ieee_method_1(double x) {
43     char * result = new char[64];
44     int exponent = 1023;
45
46     memset(result, '0', 64);
47
48     cout << "Convert " << x << " by method 1:" << endl;
49
50     // Special case: x=0
51     // Default: x=positive zero
52     if (x == 0) {
53         return result;
54     }
55
56     // Determine sign bit
57     if (x < 0) {
58         result[0] = '1';
59         x = -x;
60     }
61
62     // Determine exponent
63     if (x >= 2) {
64         while (x >= 2) {
65             x /= 2;
66             exponent++;
```



```

67     }
68 }
69 else {
70     if (x < 1) {
71         while (x < 1) {
72             x *= 2;
73             exponent--;
74         }
75     }
76 }
77 for (int i = 11; exponent > 0; i--, exponent >= 1) {
78     result[i] = char('0' + (exponent & 1));
79 }
80
81 // Determine significant digits
82 // Since subnormal numbers are NOT considered, at this moment, x must lies in [1,2).
83 x--;
84 for (int i = 12; x > 0 && i < 64; i++) {
85     x *= 2;
86     result[i] = x >= 1 ? '1' : '0';
87     if (x >= 1) {
88         x--;
89     }
90 }
91 return result;
92 }
93
94 // Convert by definition (IEEE 754)
95 // ** Subnormal numbers are NOT considered **
96 double ieee_2_double_method_1(char * s) {
97     double result = 1;
98     double temp = 1;
99     int exponent = 0;
100
101     cout << "Convert back by method 1:" << endl;
102
103     // Determine significant digits
104     for (int i = 12; i < 64; i++) {
105         temp /= 2;
106         result += temp * (int(s[i]) - int('0'));
107     }
108
109     // Determine exponent
110     for (int i = 11; i > 0; i--) {
111         exponent += (int(s[i]) - int('0')) * (1 << (11 - i));
112     }
113     if (exponent == 0) {
114         return 0;
115     }
116     exponent -= 1023;
117     if (exponent >= 0) {
118         for (; exponent > 0; exponent--) {
119             result *= 2;
120         }
121     }
122     else {
123         for (; exponent < 0; exponent++) {
124             result /= 2;
125         }
126     }
127
128     // Determine sign
129     if (s[0] == '1') {
130         result = -result;
131     }
132
133     return result;
134 }
135
136 // Convert by explicit type conversion
137 char * double_2_ieee_method_2(double x) {
138     char * result = new char[64];
139     long long * byte = (long long *)(&x);    // sizeof(long long)=sizeof(double)=8

```

```

140
141     memset(result, '0', 64);
142
143     cout << "Convert " << x << " by method 2;" << endl;
144
145     // Fill result from the back 'bit' by 'bit'
146     for (int i = 63; i >= 0; i--) {
147         result[i] = (char)(int('0') + ((*byte) & 1));
148         *byte >>= 1;
149     }
150
151     return result;
152 }
153
154 // Convert by explicit type conversion
155 double ieee_2_double_method_2(char * s) {
156     long long result = 0;
157
158     cout << "Convert back by method 2:" << endl;
159
160     // Convert the string to long long
161     for (int i = 0; i < 64; i++) {
162         result <<= 1;
163         result += (int)(s[i]) - (int)'0';
164     }
165
166     // Converting the long long to a double
167     return *((double *)&result);
168 }
169
170 // Method for output of cstrings
171 inline void stringout(char * s, int n) {
172     for (int i = 0; i < n; i++) {
173         cout << s[i];
174     }
175     cout << endl;
176 }
177
178 // Method for showcasing the examples
179 inline void showcase(double x) {
180     char * s1, *s2;
181     stringout(s1 = double_2_ieee_method_1(x), 64);
182     stringout(s2 = double_2_ieee_method_2(x), 64);
183     cout << setprecision(15) << ieee_2_double_method_1(s1) << endl;
184     cout << setprecision(15) << ieee_2_double_method_2(s2) << endl;
185 }
186
187 int main() {
188     double x = 0;
189     char * s = new char[64];
190     // Example from Chapter 1 Lecture note Page 50
191     cout << "Example from Chapter 1 Lecture note Page 50" << endl;
192     showcase(-222.111);
193
194     // Other examples
195     cout << "Try another example. Input the double:" << endl;
196     cin >> x;
197     showcase(x);
198
199     cout << "Try another example. Input the string:" << endl;
200     cin >> s;
201     cout << setprecision(15) << ieee_2_double_method_1(s) << endl;
202     cout << setprecision(15) << ieee_2_double_method_2(s) << endl;
203     system("pause");
204     return 0;
205 }

```

A.1.2 第三题

```

1 #include <iostream>
2 #include <iomanip>
3 #include <fstream>
4
5 using namespace std;

```

```

6
7 double a, b, tau_a, tau_b, t, delta;
8
9 double a_prev;
10 double b_prev;
11 ofstream out("result.txt");
12
13 void solve_core() {
14 // Method 1
15 // Naive iteration
16 /*
17     a -= delta * a_prev / tau_a;
18     b += delta * (a_prev / tau_a - b_prev / tau_b);
19 */
20 // Method 2
21 // Naive iteration for a, a small correction for b
22 /*
23     a -= delta * a_prev / tau_a;
24     b += delta * ((a_prev + a) / (2 * tau_a) - b_prev / tau_b);
25 */
26 // Method 3
27 // RK4 for a, no correction for b
28 /*
29     double k[4];
30     k[0] = - delta * a_prev / tau_a;
31     k[1] = - delta * (a_prev + k[0] / 2) / tau_a;
32     k[2] = - delta * (a_prev + k[1] / 2) / tau_a;
33     k[3] = - delta * (a_prev + k[2]) / tau_a;
34     a += (k[0] + 2 * k[1] + 2 * k[2] + k[3]) / 6;
35     b += delta * ((a_prev + a) / (2 * tau_a) - b_prev / tau_b);
36 */
37 // Method 4
38 // RK4 for a, after that, RK4 for b
39
40     double k[4];
41     k[0] = -delta * a_prev / tau_a;
42     k[1] = -delta * (a_prev + k[0] / 2) / tau_a;
43     k[2] = -delta * (a_prev + k[1] / 2) / tau_a;
44     k[3] = -delta * (a_prev + k[2]) / tau_a;
45     a += (k[0] + 2 * k[1] + 2 * k[2] + k[3]) / 6;
46
47     k[0] = delta * (a_prev / tau_a - b_prev / tau_b);
48     k[1] = delta * ((a_prev + a) / (2 * tau_a) - (b_prev + k[0] / 2) / tau_b); // Explicit
49     k[2] = delta * ((a_prev + a) / (2 * tau_a) - (b_prev + k[1] / 2) / tau_b); // Explicit
50     k[3] = delta * (a / tau_a - (b_prev + k[2]) / tau_b);
51     b += (k[0] + 2 * k[1] + 2 * k[2] + k[3]) / 6;
52 }
53 }
54
55 void solver() {
56     double t_now = 0;
57     int n = 0;
58     cin >> t;
59     cin >> delta;
60     out << "a = " << setw(16) << setprecision(4) << a << endl;
61     out << "b = " << setw(16) << setprecision(4) << b << endl;
62     out << "tau_a = " << setw(16) << setprecision(4) << tau_a << endl;
63     out << "tau_b = " << setw(16) << setprecision(4) << tau_b << endl;
64     out << "t = " << setw(16) << setprecision(4) << t << endl;
65     out << "delta = " << setw(16) << setprecision(4) << delta << endl;
66     a_prev = a;
67     b_prev = b;
68     while (t_now <= t) {
69         n++;
70         out << setw(16) << setprecision(4) << t_now << setw(16) << setprecision(4) << a <<
71         setw(16) << setprecision(4) << b << endl;
72         solve_core();
73         t_now += delta;
74         a_prev = a;
75         b_prev = b;
76     }
77 }

```

```

76     out << setw(16) << setprecision(6) << t_now << setw(16) << setprecision(4) << a << setw
    (16) << setprecision(4) << b << endl;
77     a = ((t - t_now) / delta + 1) * (a - a_prev) + a;
78     b = ((t - t_now) / delta + 1) * (b - b_prev) + b;
79 }
80
81 int main() {
82     a = b = 1;
83     cin >> tau_a >> tau_b;
84     solver();
85     cout << setw(16) << setprecision(4) << a << setw(16) << setprecision(4) << b << endl;
86     getchar();
87     getchar();
88     return 0;
89 }

```