

上机题第一题实验报告

王伟任

计 33

2013011333

一、题目要求及分析

第一章上机题 3: 编程观察无穷级数 $\sum_{n=1}^{\infty} \frac{1}{n}$ 的求和计算。分别采用 IEEE 单精度浮点数和 IEEE 双精度浮点数计算，观察当 n 为何值时求和结果不再变化。

由题可知，当求和结果不再变化时，说明此时新增项 $\frac{1}{n}$ 的精度已经超出浮点数求和的计算精度。浮点数求和时，需要将两个相加的浮点数按位对齐，即小数点对齐。在这种情况下，当一个加数很大另一个加数很小时，由于浮点数所能表示的精度有限，因此求和运算会优先保留较大数的有效位数，如果较小数的有效位数对齐到浮点数精度范围之外就会被视为 0。

IEEE 单精度浮点数能够表示 24 位二进制有效位，即 24 位尾数；IEEE 双精度浮点数能够表示 53 位尾数。

二、实验结果及分析

(1) 采用 IEEE 单精度浮点数计算：

由于单精度浮点数所能表示的精度并不算大，因此直接采用依次相加判断结果是否与上一次的结果相同的办法。最后输出的结果如下：

Answer is 2097152

The summary is 15.403683

Time cost: 15

当 $n=2097152$ 时，结果中再加 $\frac{1}{n}$ 就没有效果了。此时总和约为 15.4，输出结果只保留六位小数。得出这些结果花费 15 个计时单位，即 15ms。

2097152 是 2 的 21 次方，因此 $1/2097152$ 用二进制表示为 0.00000000000000000001，小数点后共 20 个 0，1 在小数点后第 21 位。而当前的级数和整数部分为 15，用二进制表示即小数点前有四位 1111。单精度浮点数能够表示 24 位有效位数，因此求和的时候小数点后至多能够记录 20 位；然而加数 $1/2097152$ 的首位非零位在小数点后第 21 位，与已有求和结果对齐时，第 21 位由于超过记录范围被截断，因此相当于加上 0，结果没有变化。

理论上，从 2 的 20 次方开始，加数用二进制表示结果都是小数点后跟着 20 个 0，小数点后第 21 位开始出现非零位，这些数应该都被舍去。但是，C++ 编译中会将这些非零位数大于等于两位的数取近似为 0.00000000000000000001，小数点后有 19 个 0，第 20 位变成了 1。所以说，2097151 以及之前的一百多万个数在加到已有求和结果上时加的数实际都是 $9.5367431640625e-07$ ，也就是 $2^{(-20)}$ 。

(2) 采用 IEEE 双精度浮点数计算：

理论上，双精度浮点数的计算过程应该和单精度浮点数相同，但是由于结果太大，难以用现有计算机在短时间内实现，因此只能采取近似方法。

数学上，调和级数求和有近似公式 $\sum_{k=1}^n \frac{1}{k} = \ln(n+1) + \varepsilon$ ，其中 ε 为欧拉常数。因


```

    } while(s1 != s2);
    finish = clock();
    printf("Answer is %ld\n", i);
    printf("The summary is %f\n", s1);
    printf("Time cost: %ld\n", finish - start);
    return 0;
}

```

双精度算法实现:

```
#include <stdio>
```

```
#include <cmath>
```

```

int main(){
    double s1 = 0.0;
    long long i = 1;
    double adder;
    for (;i < 20000000; i++){
        adder = 1.0 / i;
        s1 += adder;
    }
    double klans = s1;
    double ansSum = s1;
    adder = log(5*i)-log(i);
    double s2 = s1 + adder;
    adder = 1.0 / (5*i);
    s1 = s2 + adder;
    while(s1 != s2){
        ansSum = s1;
        i = 5 * i + 1;
        adder = log(5*i)-log(i);
        s2 = s1 + adder;
        adder = 1.0 / (5*i);
        s1 = s2 + adder;
    }
    long long beginn = i;
    long long eend = 5 * i;
    while (beginn < eend){
        long long delta = (eend - beginn) / 4;
        long long b[5];
        b[0] = beginn - 1;
        b[1] = beginn + delta;
        b[2] = b[1] + delta;
        b[3] = b[2] + delta;
        b[4] = eend;
    }
}

```

```

        for (int kk = 0; kk < 4; kk++){
            adder = log(b[kk + 1]) - log(b[kk] + 1);
            double as1 = ansSum + adder;
            double as2 = as1 + 1.0 / b[kk + 1];
            if (as1 == as2){
                beginn = b[kk] + 1;
                eend = b[kk + 1];
                break;
            }
            ansSum = as2;
        }
    }
    printf("Answer is %lld\n", beginn);
    printf("The summary is %.50f\n", ansSum);
    printf("Cor summary is %.50f\n", klans + log(beginn+1) - log(20000000));
    return 0;
}

```

计算结果比较:

```
#include <stdio>
```

```

int main(){
    float s1 = 0, adder1;
    double s2 = 0, adder2;

    for (int i = 1; i <= 2097152; i++){
        adder1 = 1.0 / i;
        adder2 = 1.0 / i;
        s1 += adder1;
        s2 += adder2;
    }

    printf("Using single: %.50f\n", s1);
    printf("Using double: %.50f\n", s2);
    return 0;
}

```