

Informatik - AS19

Exercise 5: Floating Point Calculations & Basic Functions

Handout: 14. Okt. 2019 06:00

Due: 21. Okt. 2019 18:00

Task 1: Floating-Point Number Representation

[Open Task](#)

This task is a text based task. You do not need to write any program/C++ file: the answer should be written in main.md (and might include code fragments if questions ask for them).

Task

We examine the normalized binary representation of floating point numbers ($\beta = 2$). Your system has a precision $p = 4$ and an exponent range of $e \in [-3; 3]$.

1. Express this setting in the lecture notation $F^*(\beta, p, e_{min}, e_{max})$.
2. Convert the following decimal numbers to the *normalized binary representation*. For each number, choose the appropriate exponent e and round to the nearest value if you cannot represent the exact value.
 - 3.1416_{10}
 - 2.718_{10}
 - 7_{10}
 - 0.11_{10}
3. For each number of (2), convert the binary representation back to their decimal form, and determine the absolute rounding error of the conversion.
4. Calculate $2.718_{10} + 3.1416_{10} + 0.11_{10}$ in the binary representation. What do you observe?

Note: To round a floating point number use binary arithmetic rounding similar to decimal arithmetic rounding, i.e. round up for a 1 and down for a 0. For example, if we round to five significant digits:

- 11.001011_2 is rounded down because the sixth significant digit is a 0 and

therefore is truncated to 11.001_2 .

- 11.001101_2 is rounded up because the sixth significant digit is a 1 and therefore is $11.001_2 + 0.001_2 = 11.010_2$.

Task 2: Point on Parabola ?

[Open Task](#)

Task

Consider a parabola (\mathcal{P}) defined as $y = g(x)$, with $g(x) = 0.9 \cdot x^2 + 1.3 \cdot x - 0.7$.

Write a program that determines if a point (x, y) lies on parabola (\mathcal{P}) or not. The input is provided by two decimal numbers in the sequence x, y . The program must output yes, if the point lies on the parabola, otherwise no. Use datatype `double` for all variables and numbers used in the calculation of $g(x)$.

You will notice that a straight forward approach (comparing for equality) does not work, i.e., for some points that clearly should be on parabola g such an approach returns result no.

Hint: Look at the difference between the exact values of the function and the values that your program calculates. Change the program so that it works properly for all points that the submission system uses as test input *without hard-coding the points*.

Note: Outputting both yes and no will get you past the automatic grading, but doing so counts as hard-coded solution and results in zero points.

Task 3: Rounding

[Open Task](#)

Task

1. Implement the following rounding function that rounds a 64-bit floating point number (type `double`) to the nearest 32-bit integer (type `int`). You may assume that the type `double` complies with the IEEE standard 754. The function is only required to work correctly if the nearest integer is in the value range of the type `int`, otherwise, the return value of the function is undefined.

Note: Usage of library rounding functions (standard or others) is not allowed.

```
// PRE: x is roundable to a number in the value range of ty
// POST: return value is the integer nearest to x, or the one
//       away from 0 if x lies right in between two integers
int round(double x);
```

2. Write a program which inputs a number of type `double` from the user, then rounds this number using your function from (1), and then outputs the rounded number.

Task 4: Binary Expansion

[Open Task](#)

Task

Write a program that performs the binary expansion for a given decimal input number x , where $0 \leq x < 2$. Use the algorithm presented in the lecture. The program must output the first 16 digits of the number in the format: $b_0.b_1b_2b_3 \dots b_{15}$. Always print all 16 digits, even the trailing zeros. Do not normalize or round the number. You can structure your program into functions to avoid code repetition. Do not forget to annotate functions with pre- and post conditions.

Task 5: Fixing Functions

[Open Task](#)

This task is a text based task. You do not need to write any program/C++ file: the answer should be written in main.md (and might include code fragments if questions ask for them).

Task:

What are the problems (if any) with the following functions? Fix them and find appropriate [pre-](#) and [postconditions](#).

1. function `is_even`:

```
bool is_even(unsigned int i) {
    if (i % 2 == 0) return true;
}
```

2. function `invert`:

```
double invert(double x) {  
    double result;  
    if (x != 0) {  
        result = 1 / x;  
    }  
    return result;  
}
```

Hint: The C++ compiler does not protect you from certain types of errors. Therefore, even if you run a program in Code Expert, it is not guaranteed that the behaviour you observe is the “real” one. We have prepared a [program tracing handout](#) that shows how to execute a program with a pen and paper and which conditions indicate bugs in the executed program not caught by the C++ compiler.