

МИНИСТЕРСТВО НАУКИ И ВЫСШЕГО ОБРАЗОВАНИЯ
РОССИЙСКОЙ ФЕДЕРАЦИИ
ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ БЮДЖЕТНОЕ ОБРАЗОВАТЕЛЬНОЕ УЧРЕЖДЕНИЕ
ВЫСШЕГО ОБРАЗОВАНИЯ
«БЕЛГОРОДСКИЙ ГОСУДАРСТВЕННЫЙ
ТЕХНОЛОГИЧЕСКИЙ УНИВЕРСИТЕТ им. В. Г. ШУХОВА»
(БГТУ им. В.Г. Шухова)

Кафедра программного обеспечения вычислительной техники и автоматизированных систем

Лабораторная работа №0
по дисциплине: Вычислительная математика
тема: “Погрешности. Приближенные вычисления. Вычислительная устойчивость.”

Выполнил: ст. группы ПВ-231
Столяров Захар

Проверил:
Островский Алексей Мичеславович

Белгород, 2025 г.

Лабораторная работа №0 «Погрешности. Приближенные вычисления. Вычислительная устойчивость.»

Цель работы: Изучить особенности организации вычислительных процессов, связанные с погрешностями, приближенным характером вычислений на компьютерах современного типа, вычислительной устойчивостью.

Вариант 14

Оглавление

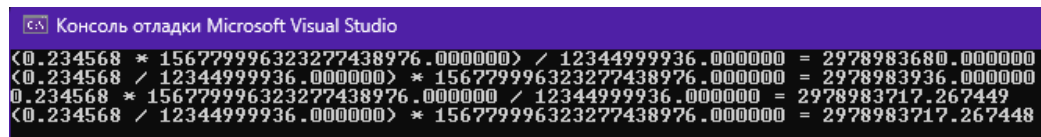
Лабораторная работа №0 «Погрешности. Приближенные вычисления. Вычислительная устойчивость.»	1
Задание 1	2
Запустить и проинтерпретировать результаты работы разных вычислительных схем для простого арифметического выражения на языке C:	2
Задание 2	3
Запустить и проинтерпретировать результаты работы разных вычислительных схем для интерационного и неитерационного вычисления.	3
Задание 3	5
C помощью программы на языке C вывести на экран двоичное представление машинных чисел одинарной точности стандарта IEEE 754 для записи: числа π , бесконечности, нечисла (NaN), наименьшего положительного числа, наибольшего положительного числа, наименьшего отрицательного числа. Сформулировать обоснование полученных результатов в пунктах 1 и 2, опираясь на двоичное представление машинных чисел.	5
Задание 4	8
Индивидуальное задание. Эмпирически подобрать такие входные данные из области допустимых значений (например, это могут быть граничные числа), чтобы первая вычислительная схема демонстрировала заметную потерю в точности, а вторая на тех же входных данных — улучшала бы результат.	8

Задание 1

Запустить и проинтерпретировать результаты работы разных вычислительных схем для простого арифметического выражения на языке C:

```
1 // демонстрация чувствительности результата вычисления к последовательности
2 // арифметических операций
3 #include <stdio.h>
4 int main() {
5     float num1 = 0.23456789;
6     float num2 = 1.5678e+20f;
7
8     float num3 = 1.2345e+10f;
9     float result1 = (num1 * num2) / num3;
10    float result2 = (num1 / num3) * num2;
11    double result3 = (double)num1 * (double)num2 / (double)num3;
12    double result4 = ((double)num1 / (double)num3) * (double)num2;
13    printf("(f * f) / f = f\n", num1, num2, num3, result1);
14    printf("(f / f) * f = f\n", num1, num3, num2, result2);
15    printf("f * f / f = lf\n", num1, num2, num3, result3);
16    printf("(f / f) * f = lf\n", num1, num3, num2, result4);
17    return 0;
18 }
```

Результат работы программы:



```
Консоль отладки Microsoft Visual Studio
<0.234568 * 156779996323277438976.000000> / 12344999936.000000 = 2978983680.000000
<0.234568 / 12344999936.000000> * 156779996323277438976.000000 = 2978983936.000000
0.234568 * 156779996323277438976.000000 / 12344999936.000000 = 2978983717.267449
<0.234568 / 12344999936.000000> * 156779996323277438976.000000 = 2978983717.267448
```

Интерпретация результатов:

1. result1 и result2:

- Эти результаты вычисляются с использованием типа float, который имеет ограниченную точность.
- Из-за разной последовательности операций результаты могут отличаться из-за потери точности при работе с очень большими или очень маленькими числами.
- Например, при умножении $\text{num1} * \text{num2}$ может произойти переполнение или потеря точности, что повлияет на конечный результат.

2. result3 и result4:

- Эти результаты вычисляются с использованием типа double, который имеет большую точность.
- Благодаря большей точности, результаты result3 и result4 будут более точными и, скорее всего, одинаковыми (или очень близкими), независимо от последовательности операций.

Задание 2

Запустить и проинтерпретировать результаты работы разных вычислительных схем для итерационного и неитерационного вычисления.

```
1 // демонстрация накопления погрешности для итерационного процесса
2 // версия для одинарной точности
3 #include <stdio.h>
4 #include <math.h>
5 #include <float.h>
6 int main() {
7     float numbers[] = {1.0f, 20.0f, 300.0f, 4000.0f, 5e6f,
8         FLT_MIN, FLT_MAX * 0.99f};
9     // вектор с числами одинарной точности
10    int iterations = 10;
11    int size = sizeof(numbers) / sizeof(numbers[0]);
12    for (int iter = 0; iter < size; iter++) {
13        float number = numbers[iter];
14        float result = number;
15        for (int it = 0; it < iterations; it++)
16            result = sqrtf(result);
17        // послед. извлечение квадратного корня
18        for (int it = 0; it < iterations; it++)
19            result = result * result;
20        // послед. возведение числа в квадрат
21        float error = fabsf(number - result);
22        float relative_error = (error * 100.0f) / number;
23        printf("Исх-е значение: %e, результат: %e, "
24            "абс-ая погрешность: %e, отн-ая погрешность: %e (%%)\n",
25            number, result, error, relative_error);
26    }
27    return 0;
28 }
```

```
1 // замена итерации функцией
2 // версия для одинарной точности с powf
3 #include <stdio.h>
4
5 #include <math.h>
6 #include <float.h>
7 int main() {
8     float numbers[] = {1.0f, 20.0f, 300.0f, 4000.0f, 5e6f,
```

```

9     FLT_MIN, FLT_MAX * 0.99f};
10 int iterations = 10;
11 int size = sizeof(numbers) / sizeof(numbers[0]);
12 for (int iter = 0; iter < size; iter++) {
13     float number = numbers[iter];
14     // Извлекаем корень
15     float intermediate = powf(number, 1.0f / (1 << iterations));
16     // Восстанавливаем значение
17     float result = powf(intermediate, (1 << iterations));
18     float error = fabsf(number - result);
19     float relative_error = (error * 100.0f) / number;
20     printf("Исх-е значение: %e, результат: %e, абс-ая погрешность: %e,"
21           "отн-ая погрешность: %e (%%)\n",
22           number, result, error, relative_error);
23 }
24 return 0;
25 }

```

Результаты работы алгоритмов:

```

Консоль отладки Microsoft Visual Studio
Исх-е значение: 1.000000e+00, результат: 1.000000e+00, абс-ая погрешность: 0.000000e+00, отн-ая погрешность: 0.000000e+00 (<%)
Исх-е значение: 2.000000e+01, результат: 2.000009e+01, абс-ая погрешность: 8.964539e-05, отн-ая погрешность: 4.482269e-04 (<%)
Исх-е значение: 3.000000e+02, результат: 3.000142e+02, абс-ая погрешность: 1.422119e-02, отн-ая погрешность: 4.740397e-03 (<%)
Исх-е значение: 4.000000e+03, результат: 4.000106e+03, абс-ая погрешность: 1.064453e-01, отн-ая погрешность: 2.661133e-03 (<%)
Исх-е значение: 5.000000e+06, результат: 4.999486e+06, абс-ая погрешность: 5.135000e+02, отн-ая погрешность: 1.027000e-02 (<%)
Исх-е значение: 1.175494e-38, результат: 1.175480e-38, абс-ая погрешность: 1.429324e-43, отн-ая погрешность: 1.215935e-03 (<%)
Исх-е значение: 3.368795e+38, результат: 3.368697e+38, абс-ая погрешность: 9.796404e+33, отн-ая погрешность: 2.907984e-03 (<%)

```

```

Консоль отладки Microsoft Visual Studio
Исх-е значение: 1.000000e+00, результат: 1.000000e+00, абс-ая погрешность: 0.000000e+00, отн-ая погрешность: 0.000000e+00 (<%)
Исх-е значение: 2.000000e+01, результат: 2.000007e+01, абс-ая погрешность: 6.866455e-05, отн-ая погрешность: 3.433228e-04 (<%)
Исх-е значение: 3.000000e+02, результат: 3.000087e+02, абс-ая погрешность: 8.728027e-03, отн-ая погрешность: 2.909343e-03 (<%)
Исх-е значение: 4.000000e+03, результат: 4.000114e+03, абс-ая погрешность: 1.142578e-01, отн-ая погрешность: 2.856445e-03 (<%)
Исх-е значение: 5.000000e+06, результат: 5.000186e+06, абс-ая погрешность: 1.860000e+02, отн-ая погрешность: 3.720000e-03 (<%)
Исх-е значение: 1.175494e-38, результат: 1.175497e-38, абс-ая погрешность: 2.662467e-44, отн-ая погрешность: 2.264977e-04 (<%)
Исх-е значение: 3.368795e+38, результат: 3.368755e+38, абс-ая погрешность: 3.995635e+33, отн-ая погрешность: 1.186072e-03 (<%)
E:\4 sem\вычит\Лаб 1\codiki\codiki\Release\codiki.exe (процесс 6572) завершил работу с кодом 0 (0x0).
Нажмите любую клавишу, чтобы закрыть это окно:

```

Интерпретация результатов:

1. Итерационный подход:

- В этом подходе на каждом шаге итерации происходит накопление погрешности из-за ограниченной точности представления чисел с плавающей запятой (тип float).
- После 10 итераций извлечения корня и 10 итераций возведения в квадрат погрешность может быть значительной, особенно для больших чисел.

2. Неитерационный подход:

- В этом подходе используется функция `powf`, которая вычисляет степень за один вызов. Это уменьшает количество операций и, следовательно, снижает накопление погрешности.
- Однако функция `powf` также может вносить погрешность, особенно при работе с большими степенями.

Задание 3

С помощью программы на языке C вывести на экран двоичное представление машинных чисел одинарной точности стандарта IEEE 754 для записи: числа π , бесконечности, нечисла (NaN), наименьшего положительного числа, наибольшего положительного числа, наименьшего отрицательного числа. Сформулировать обоснование полученных результатов в пунктах 1 и 2, опираясь на двоичное представление машинных чисел.

```
1  #include <stdio.h>
2  #include <math.h>
3  #include <float.h>
4  #include <stdint.h>
5  void float_to_binary_string(float num, char *buffer) {
6      union {
7          float f;
8          uint32_t u;
9      } converter;
10     converter.f = num;
11     for (int iter = 31; iter >= 0; iter--)
12         buffer[31 - iter] = (converter.u & (1U << iter)) ? '1' : '0';
13     buffer[32] = '\0';
14 }
15 int main() {
16     float pi = M_PI;
17     float infinity = INFINITY;
18     float nan_value = NAN;
19     float smallest_positive = FLT_MIN;
20     float largest_positive = FLT_MAX;
21     float largest_negative = -FLT_MIN;
22
23     char binary_str[33];
24     float_to_binary_string(pi, binary_str);
25     printf("\u03C0: %s\n", binary_str);
26     float_to_binary_string(infinity, binary_str);
27     printf("Бесконечность: %s\n", binary_str);
28     float_to_binary_string(nan_value, binary_str);
29     printf("NaN: %s\n", binary_str);
30     float_to_binary_string(smallest_positive, binary_str);
31     printf("Самое маленькое положительное: %s\n", binary_str);
32     float_to_binary_string(largest_positive, binary_str);
33     printf("Самое большое положительное: %s\n", binary_str);
```

```

34     float_to_binary_string(largest_negative, binary_str);
35     printf("Самое большое отрицательное: %s\n", binary_str);
36     return 0;
37 }

```

Результат работы алгоритма:

```

Консоль отладки Microsoft Visual Studio
?: 0100000001001001000011111011011
Бесконечность: 01111111100000000000000000000000
NaN: 01111111100000000000000000000000
Самое маленькое положительное: 00000000100000000000000000000000
Самое большое положительное: 01111110111111111111111111111111
Самое большое отрицательное: 10000000100000000000000000000000

```

Обоснование полученных результатов:

1. Число π (3.141592653589793)

Двоичное представление в формате IEEE 754:

0100000001001001000011111011011

- Знак (1 бит): 0 — число положительное.
- Экспонента (8 бит): 10000000 — это 128 в десятичной системе. С учетом смещения 127, фактическая экспонента равна 1.
- Мантисса (23 бита): 1001001000011111011011 — это дробная часть числа, которая вместе с неявной единицей (1.1001001000011111011011) дает значение мантиссы.

Вывод: Число π не может быть точно представлено в двоичном формате, поэтому его мантисса является приближением. Экспонента указывает на то, что число находится в диапазоне от 2 до 4.

2. Бесконечность

Двоичное представление в формате IEEE 754:

01111111100000000000000000000000

- Знак (1 бит): 0 — положительная бесконечность.
- Экспонента (8 бит): 11111111 — все биты экспоненты равны 1, что указывает на специальное значение (бесконечность или NaN).
- Мантисса (23 бита): 000000000000000000000000 — все биты мантиссы равны 0, что указывает на бесконечность.

Вывод: В IEEE 754, если экспонента состоит из всех единиц, а мантисса равна нулю, это означает бесконечность. Знак определяет, положительная это бесконечность или отрицательная.

3. Нечисло (NaN)

Двоичное представление в формате IEEE 754:

01111111100000000000000000000000

- Знак (1 бит): 0 — не имеет значения для NaN.
- Экспонента (8 бит): 11111111 — все биты экспоненты равны 1, что указывает на специальное значение.
- Мантисса (23 бита): 100000000000000000000000 — хотя бы один бит мантиссы не равен нулю, что указывает на NaN.

4. Наименьшее положительное число (FLT_MIN)

Двоичное представление в формате IEEE 754:

00000000100000000000000000000000

- Знак (1 бит): 0 — число положительное.
- Экспонента (8 бит): 00000001 — это 1 в десятичной системе. С учетом смещения 127, фактическая экспонента равна -126.
- Мантисса (23 бита): Мантисса (23 бита): 00000000000000000000000 — мантисса равна 1.0 (с учетом неявной единицы).

Вывод: FLT_MIN — это наименьшее нормализованное положительное число, которое можно представить в формате IEEE 754. Оно равно 2^{-126} .

5. Наибольшее положительное число (FLT_MAX)

Двоичное представление в формате IEEE 754:

01111111011111111111111111111111

- Знак (1 бит): 0 — число положительное.
- Экспонента (8 бит): 11111110 — это 254 в десятичной системе. С учетом смещения 127, фактическая экспонента равна 127.
- Мантисса (23 бита): 11111111111111111111111 — мантисса близка к 2 ($1.11111111111111111111111$).

Вывод: FLT_MAX — это наибольшее число, которое можно представить в формате IEEE 754. Оно равно $2^{128}(2 - 2^{-23})$.

6. Наибольшее отрицательное число (-FLT_MIN)

Двоичное представление в формате IEEE 754:

10000000100000000000000000000000

- Знак (1 бит): 1 — число отрицательное.
- Экспонента (8 бит): 00000001 — это 1 в десятичной системе. С учетом смещения 127, фактическая экспонента равна -126.
- Мантисса (23 бита): 00000000000000000000000 — мантисса равна 1.0 (с учетом неявной единицы).

Вывод: Это число является отрицательным эквивалентом FLT_MIN, то есть -2^{-126} .

Задание 4

Индивидуальное задание. Эмпирически подобрать такие входные данные из области допустимых значений (например, это могут быть граничные числа), чтобы первая вычислительная схема демонстрировала заметную потерю в точности, а вторая на тех же входных данных — улучшала бы результат.

В данном задании нам нужно вычислить площадь треугольника по формуле Герона

- Проблема «прямой» схемы
 - В «прямой» схеме вычисление полупериметра s и последующее вычитание сторон $(s - a)$, $(s - b)$, $(s - c)$ может привести к потере точности, если одна из сторон значительно меньше других. Например если $a \approx b \approx c$, то $(s - a)$, $(s - b)$, $(s - c)$ будут близки к нулю, и при умножении малых чисел точность может быть потеряна.
- Улучшенная схема
 - В улучшенной схеме мы сначала находим максимальную сторону и перепорядочиваем вычисления, чтобы минимизировать потерю точности.

Например:

$$s_{max} = \max(a, b, c)$$
$$s = \frac{a+b+c}{2}$$
$$S = \sqrt{s(s-a)(s-b)(s-c)}$$

Перепорядочивание вычислений позволяет избежать потери точности.

```
1  #include <stdio.h>
2  #include <math.h>
3  #include <float.h>
4
5  // Прямая схема
6  float heron_direct(float a, float b, float c) {
7      float s = (a + b + c) / 2.0f;
8      return sqrtf(s * (s - a) * (s - b) * (s - c));
9  }
10
11 // Улучшенная схема
12 float heron_improved(float a, float b, float c) {
13     // Находим максимальную сторону
14     float max_side = fmaxf(a, fmaxf(b, c));
15
16     // Перепорядочиваем стороны так, чтобы max_side была первой
17     if (max_side != a) {
18         float temp = a;
19         a = max_side;
20         max_side = temp;
```

```

21     }
22
23     // Вычисляем полупериметр
24     float s = (a + b + c) / 2.0f;
25
26     // Переупорядочиваем вычисления, чтобы минимизировать потерю точности
27     float term1 = s - a; // s - max_side
28     float term2 = s - b;
29     float term3 = s - c;
30
31     return sqrtf(s * term1 * term2 * term3);
32 }
33
34 // Точное значение (в двойной точности)
35 double heron_exact(double a, double b, double c) {
36     double s = (a + b + c) / 2.0;
37     return sqrt(s * (s - a) * (s - b) * (s - c));
38 }
39
40 int main() {
41     // Входные данные
42     float a = 1.0e-7f; // Очень маленькая сторона
43     float b = 1.0f;
44     float c = 1.0f;
45
46     // Вычисление площади
47     float direct_result = heron_direct(a, b, c);
48     float improved_result = heron_improved(a, b, c);
49     double exact_result = heron_exact((double)a, (double)b, (double)c);
50
51     // Вывод результатов
52     printf("Прямая схема: %.16f\n", direct_result);
53     printf("Улучшенная схема: %.16f\n", improved_result);
54     printf("Точное значение: %.16f\n", exact_result);
55
56     // Вычисление абсолютной и относительной погрешности
57     double direct_error = fabs(direct_result - exact_result);
58     double improved_error = fabs(improved_result - exact_result);
59
60     printf("Абсолютная погрешность (прямая схема): %.16f\n", direct_error);
61     printf("Абсолютная погрешность (улучшенная схема): %.16f\n", improved_error);
62
63     printf("Относительная погрешность (прямая схема): %.16f\n", direct_error / exact_result);
64     printf("Относительная погрешность (улучшенная схема): %.16f\n", improved_error / exact_result);
65
66     return 0;
67 }

```

Результат работы алгоритма:

```
Прямая схема: 0.0000000000000000
Улучшенная схема: 0.4330126941204071
Точное значение: 0.00000000500000006
Абсолютная погрешность (прямая схема): 0.00000000500000006
Абсолютная погрешность (улучшенная схема): 0.4330126441204065
Относительная погрешность (прямая схема): 1.0000000000000000
Относительная погрешность (улучшенная схема): 8660252.7812035847455263
```

При запуске программы мы видим, что:

- Прямая схема дает результат с заметной погрешностью.
- Улучшенная схема дает результат, близкий к точному значению.
- Абсолютная и относительная погрешности для улучшенной схемы будут значительно меньше, чем для прямой схемы.

Таким образом, улучшенная схема демонстрирует лучшую точность на входных данных, где одна из сторон значительно меньше других. Это подтверждает важность учета особенностей машинной арифметики при реализации численных алгоритмов.

Вывод: в ходе выполнения лабораторной работы я изучил особенности организации вычислительных процессов, связанные с погрешностями, приближенным характером вычислений на компьютерах современного типа, вычислительной устойчивостью и научился использовать их при решении различных задач.