# Floating Point Numbers

Floating-point numbers in Java are numbers that have a decimal point or are expressed using exponential notation. Java provides two primary data types for representing floating-point numbers: `float` and `double`.

## Floating-Point Representation

Floating-point numbers in computers are represented using a binary format, which is different from our everyday decimal (base-10) system. When we talk about floating-point numbers, we're referring to numbers that have a decimal (or binary) point. The IEEE 754 standard is commonly used for this representation.

**Decimal vs. Binary**

In the decimal system, numbers are based on powers of 10. For instance, the fraction $0.125$ in decimal is represented as:

$$0.125 = \frac{1}{10} + \frac{2}{100} + \frac{5}{1000}$$

In the binary system, numbers are based on powers of 2. For example, the binary fraction $0.001$ is represented as:

$$0.001 = 0 \times \frac{1}{2^1} + 0 \times \frac{1}{2^2} + 1 \times \frac{1}{2^3}$$

When evaluated, this gives:

$$0.001 = 0 + 0 + \frac{1}{8}$$

So, the binary fraction $0.001$ is equivalent to the decimal fraction $\frac{1}{8}$ or $0.125$.

## Approximations in Binary

Most decimal fractions can't be represented exactly in binary. This is due to the differences in the base systems used. As a result, when you input a decimal floating-point number, the computer stores an approximation of it in binary.

For instance, the decimal fraction $0.1$ doesn't have an exact binary representation. If we tried to represent it, we'd get a repeating binary fraction:

$$0.1 \approx 0.0001100110011001100110011001100110011001100110011...$$

This sequence $1100$ would continue indefinitely.

## Floating-Point Precision and Arithmetic

When you work with floating-point numbers in programming languages like JavaScript, Python, Rust and Java, you're dealing with approximations. This is because most decimal numbers cannot be represented exactly in binary.

## The Case of $0.1 + 0.2$

The phenomenon where $0.1 + 0.2$ does not exactly equal $0.3$ in many programming languages is a direct consequence of how floating-point numbers are represented in binary. Take the numbers $0.1$ and $0.2$ as an example. In binary, neither of these numbers has an exact representation. When you add their approximate binary representations together, the result is an approximation of $0.3$, but not exactly $0.3$.

In JavaScript:

```
console.log(0.1 + 0.2);  // Outputs: 0.30000000000000004
```

```
0.30000000000000004
```

In Python:

```
print(0.1 + 0.2)  # Outputs: 0.30000000000000004
```

```
0.30000000000000004
```

In Java:

```
class Main {
    public static void main(String[] args) {
        System.out.println(0.1 + 0.2);  // Outputs: 0.30000000000000004
    }
}
```

**Why Does This Happen?**

The root of this issue lies in the IEEE 754 standard for floating-point arithmetic, which is used by many programming languages and computer architectures. In this standard, numbers are stored in a binary format, which can't represent $0.1$ or $0.2$ (and many other numbers) exactly. When these approximations are added together, the result is slightly off from what we might expect.

## Implications

This behavior can lead to unexpected results in calculations, especially when comparing numbers for equality. For instance, checking if $0.1 + 0.2 == 0.3$ will return `false` in many languages. Developers need to be aware of this behavior and use techniques like setting a small tolerance level when comparing floating-point numbers.

## Java solution for precision issues

The `BigDecimal` class in Java is designed to address the precision issues inherent in binary floating-point representations, such as the one with `0.1 + 0.2`.

Using `BigDecimal`, you can perform arithmetic with decimal numbers and get exact results. Here's how you can use `BigDecimal` to solve the `0.1 + 0.2` issue:

```java
import java.math.BigDecimal;

public class Main {
    public static void main(String[] args) {
        BigDecimal a = new BigDecimal("0.1");
        BigDecimal b = new BigDecimal("0.2");
        BigDecimal result = a.add(b);

        System.out.println(result);  // Outputs: 0.3
    }
}
```

Note a few things about using `BigDecimal` :

1. **String Constructor**: It's recommended to use the `BigDecimal` constructor that takes a `String` argument rather than the one that takes a `double` . This is because the `double` value might already have precision issues by the time it's passed to the `BigDecimal` constructor.
2. **Immutable**: `BigDecimal` objects are immutable. Operations on them return a new `BigDecimal` instance rather than modifying the current instance.
3. **Performance**: While `BigDecimal` provides exact precision, it's slower than primitive floating-point arithmetic. It's essential to consider this trade-off, especially in performance-critical applications.

If you need exact decimal arithmetic in Java, `BigDecimal` is the way to go. It's especially useful in financial applications where precision is crucial.

# Data Types

| Data Type | Size | Range | Precision | Example |
|---|---|---|---|---|
| `float` | 32 bits (4 bytes) | Approximately ±3.4E-38 to ±3.4E+38 | About 7 decimal digits | `3.14F` |
| `double` | 64 bits (8 bytes) | Approximately ±4.9E-324 to ±1.8E+308 | About 15 decimal digits | `3.141592653589793` |

## float

The `float` data type is a single-precision 32-bit IEEE 754 floating-point. When defining a `float` literal, you use the `F` or `f` suffix to specify that the number is a float. For example: `float pi = 3.14F` .

## double

The `double` data type is a double-precision 64-bit IEEE 754 floating-point. When defining a `double` literal, you can use the `D` or `d` suffix to specify that the number is a double.

```java
class Main {
    public static void main(String[] args) {
        // Java by default treats any floating-point number as a double
        double pi = 3.141592653589793;
        System.out.println(pi);   // Outputs: 3.141592653589793

        // using the d suffix
        double e = 2.718281828459045d;
        System.out.println(e);   // Outputs: 2.718281828459045

        // using the D suffix
        double phi = 1.618033988749895D;
        System.out.println(phi);   // Outputs: 1.618033988749895

    }
}
```

> ✏️ **Note**
>
> By default, any floating-point number without a suffix is treated as a `double`.

# Scientific Notation

One of the ways to create floating-point numbers in programming is to use scientific notation. Scientific notation is a compact and convenient way to represent very large or very small numbers using a combination of a base number and an exponent. In Java, for instance, you can use the letter `e` to denote scientific notation. When you use `e`, it's typically treated as a `double` unless you explicitly specify it as a `float` by appending an 'f' or 'F'. For example, if we have the number `5e1`, it is equivalent to `50.0` in decimal form and is treated as a `double` by default. If we want to declare it as a `float`, we can write it as `5e1f` or `5e1F`. Scientific notation is useful when you want to represent very large or very small numbers compactly.

```java
class Main {
    public static void main(String[] args) {
        // 5e1 is equivalent to 50.0
        System.out.println(5e1);  // Outputs: 50.0

        // 5e1f is equivalent to 50.0f
        System.out.println(5e1f);  // Outputs: 50.0

        // 5e1F is equivalent to 50.0F
        System.out.println(5e1F);  // Outputs: 50.0

        // 5e-1 is equivalent to 0.5
        System.out.println(5e-1);  // Outputs: 0.5

        // 5e-1f is equivalent to 0.5f
        System.out.println(5e-1f);  // Outputs: 0.5

        // 5e-1F is equivalent to 0.5F
        System.out.println(5e-1F);  // Outputs: 0.5

        System.out.println(5e18);  // Outputs: 5.0E18 - This is easier to read than 500
    }
}
```

# Special Values

**Negative Infinity ( `-Infinity` ):**

- This value represents a number that's smaller than any other `float` or `double` value.
- It can result from operations like dividing a negative number by zero.
- In Java, it can be represented as `Float.NEGATIVE_INFINITY` or `Double.NEGATIVE_INFINITY` .

**Positive Infinity ( `Infinity` ):**

- Represents a number larger than any other `float` or `double` value.
- It can result from operations like dividing a positive number by zero.
- In Java, it's represented as `Float.POSITIVE_INFINITY` or `Double.POSITIVE_INFINITY` .

**Positive Zero ( `+0.0` ):**

- It's the result of certain mathematical operations, like `1.0 / Infinity` .

- While it's numerically equal to negative zero, they are distinct values in the IEEE 754 floating-point standard that Java follows.

**Negative Zero ( `-0.0` ):**

- It can result from operations like `1.0 / -Infinity` .
- Though it's numerically equal to positive zero, certain operations can distinguish between the two.

**Not-a-Number (NaN):**

- Represents a value that isn't a valid number.
- It can result from operations that don't produce a meaningful result, like taking the square root of a negative number or dividing zero by zero.
- In Java, it's represented as `Float.NaN` or `Double.NaN` .
- Any operation involving `NaN` always results in `NaN` .

**Examples**:

```java
class Main {
    public static void main(String[] args) {
        System.out.println(1.0 / 0.0);  // Outputs: Infinity
        System.out.println(-1.0 / 0.0);  // Outputs: -Infinity
        System.out.println(1.0 / Double.POSITIVE_INFINITY);  // Outputs: 0.0
        System.out.println(1.0 / Double.NEGATIVE_INFINITY);  // Outputs: -0.0
        System.out.println(0.0 / 0.0);  // Outputs: NaN

    }
}
```

# Usage

For most applications, using `double` is recommended because of its higher precision. However, if memory is a concern, or if you're certain that the added precision of `double` isn't necessary, then `float` might be more appropriate.

# Underscores in floating point Literals

In Java, you can use underscores in numeric literals (both integer and floating-point) to make them more readable, especially for long numbers. This feature was introduced in Java 7.

For floating-point literals, underscores can be placed anywhere between digits, including between the integral and fractional parts of the number. However, there are some rules to follow:

1. You cannot place underscores at the beginning or end of a number.
2. You cannot place underscores adjacent to a decimal point.
3. You cannot place underscores prior to an `F`, `f`, `D`, or `d` suffix for `float` and `double` literals, respectively.
4. You cannot place underscores in positions where a string of digits is expected.

| Data Type | Example with Underscores | Value |
|---|---|---|
| `float` | 3.141_5F | 3.1415 |
| `double` | 6.022_140_76e23 | $6.02214076 \times 10^{23}$ |

Here are some valid examples of using underscores in floating-point literals:

```java
class Main{
    public static void main(String[] args) {
        float piApproximation = 3.141_592_653_589_793F;
        double largeDouble = 123_456.789_012_345D;
        double anotherDouble = 1_2_3_4_5.6_7_8_9;

        System.out.println(piApproximation);  // Outputs: 3.1415927
        System.out.println(largeDouble);   // Outputs: 123456.789012345
        System.out.println(anotherDouble);   // Outputs: 12345.6789

    }
}
```

And here are some invalid usages:

```java
float invalid1 = _3.14F;   // Invalid: underscore at the beginning
float invalid2 = 3.14_F;   // Invalid: underscore before the suffix
double invalid3 = 123_.456;   // Invalid: underscore adjacent to the decimal point
double invalid4 = 123._456;   // Invalid: underscore adjacent to the decimal point
double invalid5 = 0b_1010; // Invalid: underscore in the binary prefix
```

# Operations on Floating-Point Numbers

Certainly! Floating-point operations are the arithmetic and mathematical operations you can perform on floating-point numbers. Here's an overview of these operations:

## Basic Arithmetic Operations

| Operation | Symbol | Description | Example |
|---|---|---|---|
| Addition | + | Adds two numbers. | `5.2 + 3.1 = 8.3` |
| Subtraction | − | Subtracts one number from another. | `5.2 − 3.1 = 2.1` |
| Multiplication | * | Multiplies two numbers. | `5.2 * 3 = 15.6` |
| Division | / | Divides one number by another. | `5.2 / 2 = 2.6` |
| Modulo | % | Returns the remainder of a division. | `5.2 % 2 = 1.2` |

## Comparison Operations

| Operation | Symbol | Description | Example |
|---|---|---|---|
| Equal to | == | Checks if two numbers are equal. | `5.2 == 5.2` returns `true` |
| Not equal to | != | Checks if two numbers are not equal. | `5.2 != 3.1` returns `true` |
| Greater than | > | Checks if one number is greater than another. | `5.2 > 3.1` returns `true` |
| Less than | < | Checks if one number is less than another. | `3.1 < 5.2` returns `true` |
| Greater than or equal to | >= | Checks if one number is greater than or equal to another. | `5.2 >= 5.2` returns `true` |
| Less than or equal to | <= | Checks if one number is less than or equal to another. | `3.1 <= 5.2` returns `true` |

# Other Operations with `Math`

The `Math` class in Java is primarily designed to work with `double` values. Most of its methods accept `double` arguments and return `double` results. However, you can still use it with `float` values by casting appropriately.

| Operation | Description | Example |
|---|---|---|
| Absolute | Returns the absolute value of a number. | `Math.abs(-5.2)` returns `5.2` |
| Power | Raises a number to the power of another. | `Math.pow(5.2, 2)` returns `27.04` |
| Square root | Returns the square root of a number. | `Math.sqrt(25)` returns `5` |
| Exponential | Raises the number `e` to the power of a number. | `Math.exp(1)` returns `e` |
| Logarithm | Returns the natural logarithm (base `e`) of a number. | `Math.log(5.2)` |
| Minimum | Returns the minimum of two numbers. | `Math.min(5.2, 3.1)` returns `3.1` |
| Maximum | Returns the maximum of two numbers. | `Math.max(5.2, 3.1)` returns `5.2` |
| Round | Rounds a number to the nearest integer. | `Math.round(5.2)` returns `5` |
| Floor | Rounds a number down to the nearest integer. | `Math.floor(5.8)` returns `5` |
| Ceiling | Rounds a number up to the nearest integer. | `Math.ceil(5.2)` returns `6` |

```
class Main {
    public static void main(String[] args) {
        System.out.println("Math.round(5.2) = " + Math.round(5.2));  // Outputs: 5
        System.out.println("Math.round(5.8) = " + Math.round(5.8));  // Outputs: 6
        System.out.println("Math.floor(5.8) = " + Math.floor(5.8));  // Outputs: 5.0
        System.out.println("Math.ceil(5.2) = " + Math.ceil(5.2));  // Outputs: 6.0
        System.out.println("Math.abs(-5.2) = " + Math.abs(-5.2));  // Outputs: 5.2
        System.out.println("Math.pow(5.2, 2) = " + Math.pow(5.2, 2));  // Outputs: 27.0
        System.out.println("Math.sqrt(25) = " + Math.sqrt(25));  // Outputs: 5.0
        System.out.println("Math.exp(1) = " + Math.exp(1));  // Outputs: 2.718281828459
        System.out.println("Math.log(5.2) = " + Math.log(5.2));  // Outputs: 1.64865862
        System.out.println("Math.min(5.2, 3.1) = " + Math.min(5.2, 3.1));  // Outputs:
        System.out.println("Math.max(5.2, 3.1) = " + Math.max(5.2, 3.1));  // Outputs:
    }
}
```

Floating-point operations are fundamental in many areas of computing, from basic arithmetic to complex scientific calculations. However, due to the binary representation of floating-point numbers, these operations might not always produce exact results, leading to rounding errors. It's essential to be aware of these potential inaccuracies, especially in precision-critical applications.

# Pitfalls

When working with floating-point numbers, there are a few pitfalls to be aware of. Let's look at some of them.

## Rounding Errors with Large Numbers

When working with very large numbers, rounding errors can accumulate, leading to inaccurate results. Here's an example that demonstrates rounding errors when working with large numbers:

```java
class Main {
    public static void main(String[] args) {
        double largeNumber = 1e18;  // 1 followed by 18 zeros
        double increment = 0.1;

        // Adding a small increment to a very large number multiple times
        for (int i = 0; i < 10; i++) {
            largeNumber += increment;
        }

        // Subtracting the original large number to see the accumulated result
        double result = largeNumber - 1e18;

        System.out.println("Expected result: " + (increment * 10));
        // Expected result: 1.0
        System.out.println("Actual result: " + result);
        // Actual result: 0.0
    }
}
```

In this example:

1. We start with a very large number `1e18` (which is `1` followed by `18` zeros) and a small increment of `0.1`.
2. We add this small increment to the large number ten times.
3. We then subtract the original large number to see the accumulated result of the ten increments.
4. Ideally, the result should be `1.0` (since `0.1` added ten times is `1.0`). However, due to rounding errors when working with such large numbers, the actual result might be slightly different from the expected `1.0`.