
Floating Point Numbers

- Representing *decimal* point numbers representation

- A **floating point number** can be written as **2 numbers**:

- An **exponent**, and
 - A **mantissa**

- **Example:**

- $314.159 = 0.314159 \times 10^3 = 314159 \times 10^{-3}$

(Remember that **multiply/dividing by 10** with **decimal number** can be accomplished by **shifting the decimal point one place to right/left**)

- We **do not** need to record the fact that the **exponent used is 10** (it is **assumed**)

Then, we can represent the number 314.159 using a **pair of numbers**:

- $314.159 = (314.159, 0) = (0.314159, 3) = (314159, -3)$

The **first number** is the **mantissa** and the **second** is the **exponent**.

- Representing *binary* floating point numbers

- **Binary floating point numbers** can also be written as an *exponent* and a *mantissa*, but now using **powers of 2** (instead of 10 for *decimal numbers*):

- **Example:**

- $1010.1011 = 0.10101011 \times 2^4 = 10101011 \times 2^{-4}$

(Similarly, **multiply/dividing by 2** with **binary number** can be accomplished by **shifting the decimal point one place to right/left**)

- Again, we **do not need** to record the fact that the **exponent used is 2** (because every number inside the computer is in binary).

We can represent the number 1010.1011 using a **pair of binary numbers**:

- $1010.1011 = (1010.1011, 0) = (0.10101011, 0000100 (4)) = (10101011, 1111100 (-4))$

The first number is the mantissa and the second is the exponent.

- **The IEEE Standard for floating point representation**

- IEEE is a standard making organization (Institute of Electric and Electronics Engineers).
- IEEE has defined a number of floating point number formats (single precision and double precision)
- The **IEEE Single Precision representation**:

- Uses 32 bits (4 bytes)

- **Format:**

```

          S EEEEEEEE MMMMMMMMMMMMMMMMMMMMMMM
Bit:    0 1      8 9                          31

```

S = sign of the mantissa
M = mantissa
E = exponent

- The **mantissa** uses the **sign/magnitude representation**:

- **s** = the **sign** of the mantissa
- **MMM...M** = the **magnitude** of the mantissa

- The **mantissa** is **normalized** so that:

- $1.0 \leq M < 2.0$
- The **leading digit 1** is **omitted**

- The **exponent EEEEEEE** uses the **excess 127 encoding scheme** for **signed numbers** that is **similar** to the **2's complement representation**.

The **excess 127 encoding** for **8 bits**:

Bit pattern	Encodes the value
00000000	-127
00000001	-126
....	
01111111	0
10000000	1
10000001	2
....	
11111111	128

- **Example:**

Given the following floating point representation:

```

01000000101000000000000000000000
= 0 1000001 0100000000000000000000

```


