

If you like, you can think of the exponent as specifying the number of places the decimal point was moved. 6.02×10^{23} has its decimal point shifted left 23 places. A negative exponent just means that the decimal point was shifted to the other way, so 6.63×10^{-34} had its exponent shifted to the right 34 places. Use positive exponents for big numbers, negative exponents for numbers very close to zero.

Scientific notation has the nice property that it is easy to use in multiplication and division. When you multiply two numbers in scientific notation, multiply the mantissas and add the exponents.

Example 37:

$$\begin{aligned} & 6.02 \times 10^{23} * 6.63 \times 10^{-34} \\ &= (6.02 * 6.63) \times 10^{23 + -34} \\ &= 39.9126 \times 10^{-11} \end{aligned}$$

To divide, divide the mantissas and subtract the exponents.

Example 38:

$$\begin{aligned} & 6.02 \times 10^{23} / 6.63 \times 10^{-34} \\ &= (6.02 / 6.63) \times 10^{23 - -34} \\ &= 0.908 \times 10^{57} \end{aligned}$$

3.2 IEEE 754 Floating Point Representation

Most computer manufacturers used to develop their own floating point representations for their own computers. Not only were they different, but also many had serious design errors. The IEEE 754 standard attempts to overcome these problems and has been adopted in most modern computers.

IEEE floating point numbers come in two sizes, four-byte single precision and eight-byte double precision numbers. The layouts for the parts of a floating point number are:

3.2.1 Single Precision (Bit Layout: sxxx xxxx xfff ffff ffff ffff ffff)

The IEEE 754 standard defines several number formats and precisions. The 32 bit format has a 1-bit sign, an 8-bit exponent with a bias of 127, and a 23-bit significand. The significand is always stored in “normalised” form with its most significant bit “1”. As this bit is always a 1, it is redundant and can be omitted from the stored number and automatically inserted in the arithmetic unit when calculations are to be done. The bits are used as — sxxx xxxx xfff ffff ffff ffff ffff where s is the sign bit, x...x are the exponent bits and f...f the significand bits. The value of a number is then

$$(-1)^{sign} * (1.0 + significand) * 2^{(exponent - 127)}$$

Sign Bit:

A sign bit of zero indicates a positive number and a sign bit of one indicates a negative number. The **mantissa** is always interpreted as a positive base-two number (unsigned). It is not a twos-complement number. If the sign bit is one, the floating-point value is negative, but the mantissa is still interpreted as a positive number that must be multiplied by -1.

Exponent Bits:

The exponent field is interpreted in one of three ways.

- An exponent of all ones indicates the floating-point number has one of the special values of plus or minus infinity, or "not a number" (NaN). NaN is the result of certain operations, such as the division of zero by zero.
- An exponent of all zeros indicates a denormalized floating-point number.
- Any other exponent indicates a normalized floating-point number.

Exponents that are neither all ones nor all zeros indicate the power of two by which to multiply the normalized mantissa. The power of two can be determined by interpreting the exponent bits as a positive number, and then subtracting a bias from the positive number. For a float, the bias is 127.

Mantissa Bits:

The mantissa contains one extra bit of precision beyond those that appear in the mantissa bits. **(1.0+mantissa)** The mantissa of a float, which occupies only 23 bits, has 24 bits of precision. The mantissa of a double, which occupies 52 bits, has 53 bits of precision. The exponent of floating-point numbers indicates whether or not the number is normalized. If the exponent is all zeros, the floating-point number is denormalized and the most significant bit of the mantissa is known to be a zero. Otherwise, the floating-point number is normalized and the most significant bit of the mantissa is known to be one.

Example 39:

$$\begin{aligned} & 0 \ 01111101 \ 101 \ 0000 \ 0000 \ 0000 \ 0000 \ 0000 \\ & \rightarrow \text{Sign bit} = 0 \\ & \rightarrow \text{Exponent bit} = 01111101 = 125 \\ & \rightarrow \text{Mantissa bit} = 1010...0 = 0.5 + 0.125 = 0.625 \\ & (-1)^0 * (1.0 + 0.625) * 2^{(125-127)} = (1.625) * 2^{-2} = \mathbf{0.40625} \end{aligned}$$

An exponent field in a float of 01111101 yields a power of two by subtracting the bias (127) from the exponent field interpreted as a positive integer (125). The power of two, therefore, is $2^{125 - 127}$, which is 2^{-2} . Mantissa (Significand) values are 101 0000 0000 0000 0000 0000. The answer is 0.625. Therefore the final answer is $(1.0 + 0.625)$ multiply two to the power of (-2) and it is equal to 0.40625.

Special Numbers:

The IEEE 754 standard has quite complicated rules on the rounding of numbers. It also has ways of representing underflowed and overflowed numbers and special error values called “Not a Number” (NaN), from cases like 0/0. Normalisation is also rather more complicated than is described here, to handle a “gradual underflow”.

NaN

An exponent of all ones with any other mantissa is interpreted to mean "not a number". The JVM always produces the same mantissa for NaN, which is all zeros except for the most significant mantissa bit that appears in the number (1 1111111 100000000000000000000000)

Infinity:

An exponent of all ones with a mantissa whose bits are all zero indicates infinity. The sign of the infinity is indicated by the sign bit.

Denormalized

An exponent of all zeros indicates the mantissa is denormalized, which means the unstated leading bit is a zero instead of a one.

Denormalized float values	
Value	Float bits (sign exponent mantissa)
Smallest positive (non-zero) float	0 00000000 000000000000000000000001
Smallest negative (non-zero) float	1 00000000 000000000000000000000001
Largest denormalized float	0 00000000 111111111111111111111111
Positive zero	0 00000000 000000000000000000000000
Negative zero	1 00000000 000000000000000000000000

Table: Single Precision (Reference only):

Range Name	S 1	E 8	M 23	Hexadecimal Range	Range	Decimal Range ^s
Quiet -NaN	1	11..11	11..11 : 10..01	FFFFFFF : FFC00001		
Indeterminate	1	11..11	10..00	FFC00000		
Signaling -NaN	1	11..11	01..11 : 00..01	FFBFFFF : FF800001		
-Infinity (Negative Overflow)	1	11..11	00..00	FF800000	$-(2^{-23}) \times 2^{127}$	$< -3.4028235677973365\text{E}+38$
Negative Normalized $-1.m \times 2^{(e-127)}$	1	11..10 : 00..01	11..11 : 00..00	FF7FFFF : 80800000	$-(2^{-23}) \times 2^{127}$: -2^{-126}	$-3.4028234663852886\text{E}+38$: $-1.1754943508222875\text{E}-38$
Negative Denormalized $-0.m \times 2^{(-126)}$	1	00..00	11..11 : 00..01	807FFFF : 80000001	$-(1 \cdot 2^{-23}) \times 2^{-126}$: -2^{-149} $(-(1+2^{-52}) \times 2^{-150})^*$	$-1.1754942106924411\text{E}-38$: $-1.4012984643248170\text{E}-45$ $(-7.0064923216240862\text{E}-46)^*$
Negative Underflow	1	00..00	00..00	80000000	-2^{-150} : < -0	$-7.0064923216240861\text{E}-46$: < -0
-0	1	00..00	00..00	80000000	-0	-0

+0	0	00..00	00..00	00000000	0	0
Positive Underflow	0	00..00	00..00	00000000	> 0 : 2^{-150}	> 0 : 7.0064923216240861E-46
Positive Denormalized $0.m \times 2^{(-126)}$	0	00..00	00..01 : 11..11	00000001 : 007FFFFF	$((1+2^{-52}) \times 2^{-150})^*$ 2^{-149} : $(1-2^{-23}) \times 2^{-126}$	$(7.0064923216240862E-46)^*$ 1.4012984643248170E-45 : 1.1754942106924411E-38
Positive Normalized $1.m \times 2^{(e-127)}$	0	00..01 : 11..10	00..00 : 11..11	00800000 : 7F7FFFFF	2^{-126} : $(2-2^{-23}) \times 2^{127}$	1.1754943508222875E-38 : 3.4028234663852886E+38
+Infinity (Positive Overflow)	0	11..11	00..00	7F800000	$> (2-2^{-23}) \times 2^{127}$	$> 3.4028235677973365E+38$
Signaling +NaN	0	11..11	00..01 : 01..11	7F800001 : 7FBFFFFF		
Quiet +NaN	0	11..11	10..00 : 11..11	7FC00000 : 7FFFFFFF		

Exercise:

What is the decimal value represented by the following numbers?

- a) 1 00000000 000000000000000000000000
b) 1 11111111 111111111111111111111111
c) 0 11111111 000000000000000000000000

3.2.2 Double precision numbers (Bit Layout: sxxx xxxx xxxx ffff ffff ffff ffff ffff ffff ffff ffff ffff ffff ffff ffff)

As described earlier, the 32-bit representation is barely adequate for serious computation; the precision is limited and “rounding” errors accumulate very quickly. Some computations lasting only a second or two can become quite meaningless. Also, the number range of $10^{\pm 38}$ is too small to handle some physical quantities, or formula involving them. The 754 standard therefore includes a 64-bit representation to overcome these problems.

In Java 32-bit quantities are of type float, while 64-bit are double. Floating point results are normally produced with type double and a cast is necessary to store into a float variable. IEEE 754 double precision uses a 52-bit significand (giving about 16 decimal digits of precision) and an 11-bit exponent with a bias of 1023 (a range of about $10^{\pm 300}$). The underlying principles are as for the 32-bit representation.

$$(-1)^{\text{sign}} * (1.0 + \text{significand}) * 2^{(\text{exponent}-1023)}$$

Example 40:

0 0111111101 1010 0000 0000 0000 0000 0000 0000 0000 0000 0000 0000 0000 0000 0000 0000
➔ Sign bit = 0
➔ Exponent bit = 0111111101 = 1021
➔ Mantissa bit = 1010...0 = 0.5 + 0.125 = 0.625
 $(-1)^0 * (1.0 + 0.625) * 2^{(1021-1023)} = (1.625) * 2^{-2} = 0.40625$

An exponent field in a float of 0111111101 yields a power of two by subtracting the bias (1023) from the exponent field interpreted as a positive integer (1021). The power of two, therefore, is 1021 - 1023, which is -2. Mantissa (Significand) values are 1010 0000 0000 0000 0000 0000 0000 0000 0000 0000 0000 0000. The answer is 0.625. Therefore the final answer is (1.0 + 0.625) multiply two to the power of (-2) and it is equal to 0.40625.

Table: Double Precision (for reference only):

Range Name	S 1	E 11	M 52	Hexadecimal Range	Range	Decimal Range [§]
Quiet -NaN	1	11..11	11..11 : 10..01	FFFFFFFFFFFFFFFF : FFF8000000000001		
Indeterminate	1	11..11	10..00	FFF8000000000000		
Signaling -NaN	1	11..11	01..11 : 00..01	FFF7FFFFFFFFFFFF : FFF0000000000001		
-Infinity (Negative Overflow)	1	11..11	00..00	FFF0000000000000	$< -(2 \cdot 2^{-52}) \times 2^{1023}$	-1.7976931348623158E+308
Negative Normalized $-1.m \times 2^{(e-1023)}$	1	11..10 : 00..01	11..11 : 00..00	FFEFFFFFFFFFFFFFFF : 8010000000000000	$-(2 \cdot 2^{-52}) \times 2^{1023}$: -2^{-1022}	-1.7976931348623157E+308 : -2.2250738585072014E-308
Negative Denormalized $-0.m \times 2^{(-1022)}$	1	00..00	11..11 : 00..01	800FFFFFFFFFFFFFFF : 8000000000000001	$-(1 \cdot 2^{-52}) \times 2^{-1022}$: -2^{-1074} $(-(1 \cdot 2^{-52}) \times 2^{-1075})^*$	-2.2250738585072010E-308 : -4.9406564584124654E-324 $(-2.4703282292062328E-324)^*$
Negative Underflow	1	00..00	00..00	8000000000000000	-2^{-1075} : < -0	-2.4703282292062327E-324 : < -0
-0	1	00..00	00..00	8000000000000000	-0	-0
+0	0	00..00	00..00	0000000000000000	0	0
Positive Underflow	0	00..00	00..00	0000000000000000	> 0 : 2^{-1075}	> 0 : 2.4703282292062327E-324
Positive Denormalized $0.m \times 2^{(-1022)}$	0	00..00	00..01 : 11..11	0000000000000001 : 000FFFFFFFFFFFFFFF	$((1 \cdot 2^{-52}) \times 2^{-1075})^*$ 2^{-1074} : $(1 \cdot 2^{-52}) \times 2^{-1022}$	$(2.4703282292062328E-324)^*$ 4.9406564584124654E-324 : 2.2250738585072010E-308
Positive Normalized $1.m \times 2^{(e-1023)}$	0	00..01 : 11..10	00..00 : 11..11	0010000000000000 : 7FEFFFFFFFFFFFFFFF	2^{-1022} : $(2 \cdot 2^{-52}) \times 2^{1023}$	2.2250738585072014E-308 : 1.7976931348623157E+308
+Infinity (Positive Overflow)	0	11..11	00..00	7FF0000000000000	$> (2 \cdot 2^{-52}) \times 2^{1023}$	1.7976931348623158E+308
Signaling +NaN	0	11..11	00..01 : 01..11	7FF0000000000001 : 7FF7FFFFFFFFFFFFFF		
Quiet +NaN	0	11..11	10..00 : 11..11	7FF8000000000000 : 7FFFFFFFFFFFFFFF		

[§] Your least significant digits may differ.

3.2.3 Converting from IEEE 754 Floating Point Representation to Decimal

Example 41:

```

4090000016 = 0100 0000 1001 0000 ... 0000
➔ Sign bit = 0
➔ Exponent bit = 100 0000 1 = 129

```

→ Mantissa bit = 001 0000 ... 0000 = 0.125
 Answer = $(-1)^0 * (1.0 + 0.125) * 2^{(129-127)} = (1.125) * 2^2 = 4.5$

Exercise:

Convert C2100000₁₆ from IEEE 754 Floating Point (Single Precision) to decimal

3.2.4 Convert from Decimal to IEEE 754 Floating Point Representation

Example 42:

(1) -1.25₁₀, write the number in binary format: -1.01₂
 → 1.01₂ is already in normalized format, so don't need to do any shifting
 → IEEE Sign bit = negative = 1
 → IEEE Significand bit => 0.01... = 0.01...0
 → Exponent = 127 = 01111111₂

Answer: 1 01111111 0100...0 = BFA00000₁₆

(2) 0.15625₁₀, write the number in binary format: 0.00101₂
 → Normalize => Shift point to the right for three places = 1.01₂ x 2⁻³
 → IEEE Sign bit = positive = 0
 → IEEE Significand bit => 0.01 = 0.010...0
 → Exponent = 127 + three places right shift = 127 + (-3) = 124 = 01111100₂

Answer: 0 01111100 0100000...0 = 3E200000₁₆

(3) 19.5₁₀, write the number in binary format: 10011.100...
 → Normalise => Shift point to the left for four places = 1.0011100... x 2⁴
 → IEEE Sign bit = positive = 0
 → IEEE Significand bit => 0.00111000 = 0.00111...0
 → Exponent = 127 + four places left shift = 127 + 4 = 131 = 10000011₂

Answer: 0 10000011 0011100000...0

Then, what is the value of 4.875₁₀?

→ 19.5₁₀ = 0 10000011 0011100000...0 from the above calculation
 → And 4.875₁₀ = 19.5 / 4 = 19.5 * 2⁻²
 → Sign Bit: unchanged = 0
 → Significand bit: unchanged = 001110000...0
 → Exponent bit = exponent of part (a) + -2 = 10000001₂ - 2₁₀ = 10000001₂

Answer is 0 10000001 00111000...0 = 409C0000₁₆

6.5₁₀, write the number in binary format: 110.100...
 → Normalise => Shift point to the left for two places = 1.10100... x 2²
 → IEEE Sign bit = positive = 0
 → IEEE Significand bit => 0.101000... = 0.1010...0
 → Exponent = 127 + two places left shift = 127 + 2 = 129 = 10000001₂

Answer: 0 10000001 10100000...0

Then, what is the value of 52₁₀?

→ 6.5₁₀ = 0 10000001 10100000...0 from the above calculation
 → And 52₁₀ = 6.5 * 8 = 6.5 * 2³
 → Sign Bit: unchanged = 0
 → Significand bit: unchanged = 10100000...0
 → Exponent bit = exponent of part (a) + 3 = 10000001₂ + 3₁₀ = 10000100₂

Answer is 0 10000100 10100000...0 = 42500000₁₆

Then, what is the value of 3.25₁₀?

→ And 3.25₁₀ = 6.5 / 2 = 6.5 * 2⁻¹
 → Sign Bit: unchanged = 0
 → Significand bit: unchanged = 10100000...0
 → Exponent bit = exponent of part (a) + (-1) = 10000001₂ - 1₁₀ = 10000000₂

Answer is 0 10000000 10100000...0 = 40500000₁₆

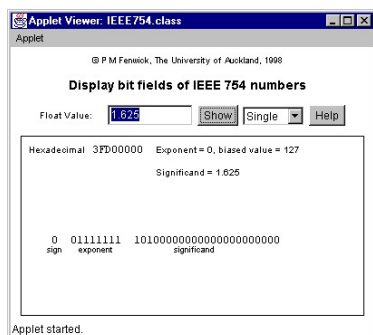
Exercises:

Convert from Decimal to IEEE 754 Floating Point (Single Precision)

(1) 2.25

(2) 4.5

3.2.5 Applet



3.2.6 IEEE 754 Floating Point in JAVA

Java provides both single (32-bit) and double (64-bit) floating point types, with the default being double.

Example 43:

```
public class IEEE754TestFloat {
    public static void main (String args[]) {
        float f = 4.5F;
        System.out.println(Integer.toHexString(Float.floatToIntBits (f)));

        f = Float.intBitsToFloat (0xbfa00000);
        System.out.println(f);
    }
}
```

Output:
> java IEEE754TestFloat
40900000
-1.25

Methods:

public static int floatToIntBits (float value)	Returns a representation of the specified floating-point value according to the IEEE 754 floating-point "single format" bit layout.
public static int floatToRawIntBits (float value)	Returns a representation of the specified floating-point value according to the IEEE 754 floating-point "single format" bit layout, preserving Not-a-Number (NaN) values.

Example 44:

```
public class IEEE754TestDouble {
    public static void main (String args[]) {
        double d = 4.5;
        System.out.println(Long.toHexString(Double.doubleToLongBits (d)));

        d = Double.longBitsToDouble (0x4012000000000000L);
        System.out.println(d);
    }
}
```

Output:
> java IEEE754TestDouble
4012000000000000
4.5

Methods:

public static long doubleToLongBits (double value)	Returns the actual bit pattern of the IEEE 754 representation of the value. (It is useful for analysing or “dismantling” a floating point number.)
public static double longBitsToDouble (long bits)	takes the bit pattern in bits and returns it as a double value. (It is useful for building a floating point number.)

3.2.7 Extra Note:

Internally, most computers use a base of 2 (i.e. the fraction is multiplied by 2 *exponent*), or less often 16 or 8. The significand is usually held in sign & magnitude form, with the exponent in say excess 127. A 0.0 value is an all-0 word. Important points to remember are

- Do not confuse the *range* and the *precision* of floating point numbers.
- The *range* is determined by the exponent and determines how close to zero or far from zero a number may be. It is closely connected to the exponent form of scientific notation. An 8-bit signed exponent can have values from -128 to +127 (-126 to +127 in IEEE 754 floating point). The smallest representable number will be about 2^{-128} and the largest 2^{+127} . Remembering that $\log_2 10$ is close to 1/0.3 (page 5 of booklet), the number range is about 10^{-38} and the largest 10^{+38} . It is shown later that this range is quite inadequate for some calculations.
- The *precision* is governed by the significand (or fraction or mantissa) and gives the accuracy with which a number may be represented. Remember that N bits equals about $0.3*N$ decimal digits. A standard “32-bit real” has 23-bit precision, or not quite 7 decimal digits. A 64-bit “double” has 52-bit or 15 decimal digits. Even a 32 bit real can handle the accuracy of most physical measurements, but much of the precision is lost by rounding in lengthy calculations; this is the real justification for using 64- bit or 128-bit floating point numbers.

Floating point arithmetic is subject to rounding and truncation errors. The significand can represent only so many bits; any less significant bits must be discarded. Often, if the first discarded bit is a 1, we add 1 onto the significand to “round” the result. Thus 1.7 would round to 2, which is probably a better result than 1 (from just forgetting the bits).

Care is needed when using real-number arithmetic. Some of the problems seem to disappear with “long” numbers, but really stay there and are never more than reduced.

- The 32 bit floating point “real” on many computers is quite limited in comparison with many scientific calculators. Its range is about $10^{\pm 38}$, and its precision is not quite 7 decimal digits. Even short calculation sequences can overwhelm it.
- Arithmetic with floating point numbers is seldom exact and great care must be taken in long calculation sequences as “round-off” errors accumulate. For example a solution of a set of 40 simultaneous equations had the 3–4 least significant decimal digits quite meaningless.
- Beware of mathematical techniques which involve differences of large quantities. This is related to the previous point. Say we have two values close to 1000, both with the last decimal digit uncertain, such as $102x$ and $99x$ (both known to about 10 parts in 1000, or 1%). Subtracting gives a value $3x$, where the last digit is still uncertain, but the error is now 10 parts in 30, or about 30%. Two moderately accurate values have combined to give a value which is nearly meaningless. Some types of statistical calculation are especially sensitive to this problem.
- The result is that floating point arithmetic is not exact. Because of possible disappearance of low-order bits we cannot guarantee that $(A+B)+C = A+(B+C)$. In most cases it is very nearly true, if we are careful, but in extreme cases it is anything but true.
- Be very careful if using floating point arithmetic for financial calculations. Rounding errors may make it almost impossible to achieve reliable balances, especially if the number precision is barely adequate to represent the whole amount.