

On Reals and Floats

The set of normalized IEEE floating point numbers is discrete and finite. It has a largest element we'll call *maxfloat*, and a smallest (positive) element, *minfloat*. All of the real numbers in the interval [*minfloat*, *maxfloat*], together with their negatives, is called the floating point range (FPR):

$$\text{FPR} = \{x \in \mathbb{R} : \text{minfloat} \leq |x| \leq \text{maxfloat}\}.$$

Notice that the FPR is the (disjoint) union of two real intervals. If $|x| < \text{minfloat}$, it is said to *underflow*, while if $|x| > \text{maxfloat}$, it is said to *overflow*.

Each real number in FPR gets represented by (“rounded to”) the floating point number nearest it. Every normalized floating point number (normal float) has the form

$$f = \pm(0.b_1b_2 \dots b_t)_2 \times 2^e, \text{ where}$$

$(0.b_1b_2 \dots b_t)_2$ is interpreted as a base-2 fraction, $b_1 \neq 0$, and e is an integer satisfying $\text{minexp} \leq e \leq \text{maxexp}$. The difference between $|f|$ and the next-larger float is 2^{e-t} , and therefore the relative difference between two neighboring floats is bounded above by 2^{1-t} . The maximum relative distance between any $x \in \text{FPR}$ and its floating point representative is called the *machine epsilon*: $\epsilon_{\text{mach}} = 2^{1-t}$. The *unit round-off* is half of this: $\mu = \epsilon_{\text{mach}}/2 = 2^{-t}$.

There are also non-normal floats. These are floating point “numbers” which do not fit the scheme above. These include ± 0 , $\pm \text{Inf}$, NaN, and the subnormals (which are evenly spaced in the interval $(-\text{minfloat}, \text{minfloat})$). Most compilers allow the programmer to use subnormals, but it’s typically not the default. The common default is to set underflow to ± 0 , and overflow to $\pm \text{Inf}$, although most compilers allow the programmer to set error flags in these cases.

Arithmetic with Inf can make sense: $\text{Inf} + \text{Inf} = \text{Inf}$, $-\text{Inf} - \text{Inf} = -\text{Inf}$, $\text{PositiveNumber} * \text{Inf} = \text{Inf}$, $\text{NegativeNumber} * \text{Inf} = -\text{Inf}$, ect., but not always: NaN, which stands for “not a number”, is used for undefined results, like $0/0$, $\text{Inf} - \text{Inf}$, Inf/Inf , $0 * \text{Inf}$, etc.

With an 8-byte word (double precision), we are able to represent $2^{64} \approx 1.8 \times 10^{19}$ floats. Of these, 2^{54} are non-normal, and half of these non-normals are subnormals. Thus about 99.9% of the 8-byte floats are normal. Here $\text{maxfloat} \approx 10^{308}$, $\text{minfloat} \approx 10^{-308}$, and $\mu \approx 10^{-16}$.

With a 4-byte word (single precision), we have $2^{32} \approx 4.3 \times 10^9$ floats, about 99.2% of which are normal. Here $\text{maxfloat} \approx 10^{38}$, $\text{minfloat} \approx 10^{-38}$, and $\mu \approx 10^{-7}$.

“Number” x	Classification	$\text{fl}(x)$	Other
$\text{minfloat} \leq x \leq \text{maxfloat}$	Normal	$\pm(0.b_1b_2 \dots b_t)_2 \times 2^e$	$\text{fl}(x) = x(1 + \delta)$, $ \delta \leq \mu$
$ x < \text{minfloat}$	Underflow	± 0	subnormals...
$ x > \text{maxfloat}$	Overflow	$\pm \text{Inf}$	error flags...
$0/0$, $0 * \text{Inf}$, $\text{Inf} - \text{Inf}$, etc.	Undefined	NaN	error flags...