

Affinely extended reals with signed zero and NaN

Bill Zorn

August 29, 2017

Abstract

We present a theory for the ideal arithmetic that underlies IEEE 754 floating point. This theory is based on “typical” real arithmetic over \mathbb{R} ; it is affinely extended to include $-\infty$ and ∞ ; it represents a recoverable sign for all numbers, including 0; and it adds a special member NaN to represent values that are not a real number in the typical sense. It serves as a bridge between real arithmetic in a mathematical sense and specifications and implementations of floating point formats, particularly IEEE 754. The theory can also specify ideal, infinite-precision results of floating point computations.

Work In Progress:
This vessel is not yet seaworthy.

1 Notation

We denote the set of mathematical real numbers as \mathbb{R} . The affinely extended real numbers with signed zero and NaN, the set $\mathbb{R} \cup \{-\infty, \infty, -0, \text{NaN}\}$, we denote \mathbb{R}^* .

Numbers in \mathbb{R} we write as simple variables, a, b, x, y , while numbers in \mathbb{R}^* we distinguish with an overline, $\bar{a}, \bar{b}, \bar{x}, \bar{y}$. A written a cannot be -0 or ∞ , while a written \bar{x} can.

We distinguish arithmetic operations in the same way we distinguish numbers. $a + b$ denotes typical addition of real numbers, while $\bar{x} \bar{+} \bar{y}$ denotes addition under the arithmetic of extended real numbers with signed zero and NaN, which we define in this document. $a \bar{+} y$ is well-defined (and hopefully not too different from typical addition), but $\bar{x} + \bar{y}$ is not, as the operands might not be real numbers.

2 Ordering

\mathbb{R}^* is only partially ordered. For all numbers in \mathbb{R}^* that are also in \mathbb{R} , the typical rules for comparison of real numbers apply. The following rules define

the ordering of the special members $-\infty$, ∞ , and -0 :

$$-\infty < \bar{x} \in \mathbb{R}^* \setminus \{-\infty, \text{NaN}\} \quad (1)$$

$$\bar{x} \in \mathbb{R}^* \setminus \{\infty, \text{NaN}\} < \infty \quad (2)$$

$$-0 = 0 \quad (3)$$

$$a < 0 \iff a < -0 \quad (4)$$

$$0 < a \iff -0 < a \quad (5)$$

Equality of members in $\mathbb{R}^* \setminus \{\text{NaN}\}$ is reflexive, symmetric, and transitive, as is to be expected, and the comparisons $<$ and $>$ follow the usual symmetries. $-\infty = -\infty$, $a \in \mathbb{R} > 0 \iff a > -0$, and so on.

NaN is unordered with respect to all other elements of \mathbb{R}^* . The following rules define comparisons with it:

$$\text{NaN} \not< \bar{x} \quad (6)$$

$$\text{NaN} \not> \bar{x} \quad (7)$$

$$\text{NaN} \neq \bar{x} \quad (8)$$

Note that equality is not reflexive ($\text{NaN} \neq \text{NaN}$), and that the usual symmetry of $<$ and $>$ is broken ($\bar{x} \not< \text{NaN} \not\iff \bar{x} > \text{NaN}$).

Because \mathbb{R}^* is only partially ordered, it is useful to define another notion of identicality, $\bar{x} \equiv \bar{y}$, which captures when two members of \mathbb{R}^* are actually the same element. $\bar{x} \equiv \bar{x}$ with no exceptions, even when $\bar{x} \equiv \text{NaN}$, and $0 \not\equiv -0$ even though $0 = -0$.

3 Primitives: negation, absolute value, and sign

All numbers in \mathbb{R}^* have a recoverable sign, except for NaN, which can be thought of as not a number. Unary negation flips this sign, as follows:

$$\text{neg}(-\infty) \equiv \infty \quad (9)$$

$$\text{neg}(\infty) \equiv -\infty \quad (10)$$

$$\text{neg}(-0) \equiv 0 \quad (11)$$

$$\text{neg}(0) \equiv -0 \quad (12)$$

$$\text{neg}(a \in \mathbb{R} \setminus \{0\}) \equiv -a \quad (13)$$

$$\text{neg}(\text{NaN}) \equiv \text{NaN} \quad (14)$$

Absolute value forces this sign to be positive:

$$\text{abs}(-\infty) \equiv \infty \quad (15)$$

$$\text{abs}(\infty) \equiv \infty \quad (16)$$

$$\text{abs}(-0) \equiv 0 \quad (17)$$

$$\text{abs}(a) \equiv |a| \quad (18)$$

$$\text{abs}(\text{NaN}) \equiv \text{NaN} \quad (19)$$

We also define a sign operation to recover the sign, multiplied by 1:

$$\text{sign}(-\infty) = -1 \quad (20)$$

$$\text{sign}(\infty) = 1 \quad (21)$$

$$\text{sign}(-0) = -1 \quad (22)$$

$$\text{sign}(a) = \begin{cases} -1, & a < 0 \\ 1, & a = 0 \vee 0 < a \end{cases} \quad (23)$$

$$\text{sign}(\text{NaN}) = -1 \sqcup 1 \quad (24)$$

$-1 \sqcup 1$ denotes a nondeterministic choice between -1 and 1. The sign of NaN is not undefined: it must be -1 or 1, though this theory does not specify which one of those it is.

Finally, we can transfer the sign from one number to another:

$$\text{copysign}(\bar{x}, \bar{y}) \equiv \bar{z} \iff (\text{abs}(\bar{x}) \equiv \text{abs}(\bar{z}) \wedge \text{sign}(\bar{y}) = \text{sign}(\bar{z})) \quad (25)$$

This definition runs into some trouble with NaN. To be faithful to the IEEE 754 standard, Equation 25 must always be satisfied, even for NaN, so the non-deterministic choice must always be made in a coherent way. This could be accomplished simply by choosing a single sign for NaN (say $\text{sign}(\text{NaN}) = 1$), or by tracking each NaN in a computation and assigning some specific sign to it.

4 Addition and subtraction

All arithmetic operations in \mathbb{R}^* are closed, with a defined result in \mathbb{R}^* . Addition is symmetric:

$$\bar{x} \bar{+} \bar{y} \equiv \bar{y} \bar{+} \bar{x} \quad (26)$$

We take advantage of this to simplify the rules somewhat:

$$\text{NaN} \bar{+} \bar{x} \equiv \text{NaN} \quad (27)$$

$$\infty \bar{+} -\infty \equiv \text{NaN} \quad (28)$$

$$\infty \bar{+} \bar{x} \in \mathbb{R}^* \setminus \{-\infty, \text{NaN}\} \equiv \infty \quad (29)$$

$$-\infty \bar{+} \bar{x} \in \mathbb{R}^* \setminus \{\infty, \text{NaN}\} \equiv -\infty \quad (30)$$

$$-0 \bar{+} \bar{x} \in \mathbb{R}^* \setminus \{0\} \equiv \bar{x} \quad (31)$$

$$-0 \bar{+} 0 \equiv \begin{cases} -0, \text{roundTowardNegative} \\ 0, \neg \text{roundTowardNegative} \end{cases} \quad (32)$$

$$a \bar{+} b \equiv \begin{cases} a + b, & a + b \neq 0 \\ -0, & a + b = 0 \wedge \text{roundTowardNegative} \\ 0, & a + b = 0 \wedge \neg \text{roundTowardNegative} \end{cases} \quad (33)$$

roundTowardNegative is some arbitrary predicate representing the IEEE 754 rounding mode. If it is True, then addition producing 0 should produce -0 instead, in accordance with the IEEE specification for arithmetic with that rounding mode. It can be thought of as a parameter of this theory, or as a distinction between two similar theories.

Subtraction is defined in terms of addition and negation:

$$\bar{x} - \bar{y} \equiv \bar{x} + \text{neg}(\bar{y}) \quad (34)$$

5 Multiplication

Like addition, multiplication over \mathbb{R}^* is symmetric:

$$\bar{x} \times \bar{y} \equiv \bar{y} \times \bar{x} \quad (35)$$

It is defined as follows:

$$\text{NaN} \times \bar{x} \equiv \text{NaN} \quad (36)$$

$$\bar{x} \in \{-0, 0\} \times \bar{y} \in \{-\infty, \infty\} \equiv \text{NaN} \quad (37)$$

$$\bar{x} \in \{-\infty, \infty\} \times \bar{y} \in \mathbb{R}^* \setminus \{-0, 0, \text{NaN}\} \equiv \text{copysign}(\infty, \text{sign}(\bar{x}) \times \text{sign}(\bar{y})) \quad (38)$$

$$\bar{x} \in \{-0, 0\} \times \bar{y} \in \mathbb{R}^* \setminus \{-\infty, \infty, \text{NaN}\} \equiv \text{copysign}(0, \text{sign}(\bar{x}) \times \text{sign}(\bar{y})) \quad (39)$$

$$a \times b \equiv a \times b \quad (40)$$

Where it produces an answer other than NaN, multiplication effectively computes the sign and magnitude of the result separately; hence the use of *copysign*. This dependence could be removed by rewriting Equations 38 and 39 as cases on the signs of the operands, though this would be much more verbose.

6 Division

Division is similar to multiplication, but it cannot be defined in terms of it. We cannot simply state some equality like $\bar{x} \div \bar{y} \equiv \bar{x} \times \text{reciprocal}(\bar{y})$. Instead, we

define it as follows:

$$\text{NaN} \overline{\div} \bar{x} \equiv \text{NaN} \quad (41)$$

$$\bar{x} \overline{\div} \text{NaN} \equiv \text{NaN} \quad (42)$$

$$\bar{x} \in \{-\infty, \infty\} \overline{\div} \bar{y} \in \{-\infty, \infty\} \equiv \text{NaN} \quad (43)$$

$$\bar{x} \in \{-0, 0\} \overline{\div} \bar{y} \in \{-0, 0\} \equiv \text{NaN} \quad (44)$$

$$\bar{x} \in \{-\infty, \infty\} \overline{\div} \bar{y} \in \mathbb{R}^* \setminus \{-\infty, \infty, \text{NaN}\} \equiv \text{copysign}(\infty, \text{sign}(\bar{x}) \times \text{sign}(\bar{y})) \quad (45)$$

$$\bar{x} \in \mathbb{R}^* \setminus \{-\infty, \infty, \text{NaN}\} \overline{\div} \bar{y} \in \{-\infty, \infty\} \equiv \text{copysign}(0, \text{sign}(\bar{x}) \times \text{sign}(\bar{y})) \quad (46)$$

$$\bar{x} \in \{-0, 0\} \overline{\div} \bar{y} \in \mathbb{R}^* \setminus \{-0, 0, \text{NaN}\} \equiv \text{copysign}(0, \text{sign}(\bar{x}) \times \text{sign}(\bar{y})) \quad (47)$$

$$\bar{x} \in \mathbb{R}^* \setminus \{-0, 0, \text{NaN}\} \overline{\div} \bar{y} \in \{-0, 0\} \equiv \text{copysign}(\infty, \text{sign}(\bar{x}) \times \text{sign}(\bar{y})) \quad (48)$$

$$a \overline{\div} b \equiv a \div b \quad (49)$$

Unsurprisingly, division over \mathbb{R}^* is not symmetric. Independent sign and magnitude computations occur, as with multiplication, in Equations 45 - 48.

Division can be used to compute something like the reciprocal, but there isn't a closed reciprocal in \mathbb{R}^* : $1 \overline{\div} \infty \equiv 0$, but $\infty \overline{\div} 0 \equiv \text{NaN}$. In fact, it's impossible to find any $\bar{x} \in \mathbb{R}^*$ such that $\infty \overline{\div} \bar{x} \equiv 1$ or $0 \overline{\div} \bar{x} \equiv 1$.

7 Integer powers

An integer power function is necessary to compute the real values of floating point numbers. The implementation should be equivalent to iterated multiplication (and optionally division for negative powers).

More details coming soon!

8 Integer roots

An integer root function is useful to complete the mandated “single-ulp” correctly rounded operations in the IEEE 754 standard, $+$ $-$ $*$ $/$ and sqrt . This has not yet been implemented for Titanic, as it is not needed to convert between floating point and real numbers.

More details coming soon!