

Injectivity of fraction to float conversion

Mark Dickinson

August 20, 2022

In some situations it's convenient to store fractions as double-precision IEEE 754 binary floats. It's clear that any map from the infinite set of all fractions to the finite collection of all double-precision floats cannot be injective, but with a suitable bound on the numerator and denominator of the fractions considered, we should be able to recover injectivity of the map which takes each fraction to the closest representable float.

This note establishes that 67114657 is such a bound. This is best possible, since the distinct fractions 67114658/67114657 and 67114657/67114656 round to the same IEEE 754 binary64 float under round-to-nearest.

Theorem 1. *Suppose that a/b and c/d are rational numbers, written in lowest terms (with b and d positive), that $\max(a, b, c, d) \leq 67114657$, and that a/b and c/d become equal when rounded to the nearest finite IEEE binary64 floating-point number. Then $a/b = c/d$.*

Proof. It's straightforward to check that when we round a fraction with numerator and denominator bounded by 67114657 to a float, no underflow or overflow occurs, and that under the same bounds positive fractions round to positive floats and negative fractions to negative floats. So if a/b and c/d round to the same float then both are positive, both are negative, or both are zero, and in the last case we're done. Negating both fractions if necessary, without loss of generality we can assume going forward that all of a , b , c and d are positive.

Now we separate into two cases. The first case we consider is the case in which a/b and c/d belong to the same closed binade—that is, there's an integer e such that both a/b and c/d lie in the interval $[2^e, 2^{e+1}]$. Since consecutive IEEE 754 binary64 floats in that interval have difference 2^{e-52} , for a/b and c/d to round to the same float we must have $|a/b - c/d| \leq 2^{e-52}$, from which

$$2^{52-e}|ad - bc| \leq bd. \quad (1)$$

From this inequality combined with $2^e \leq a/b$ and $2^e \leq c/d$, we have:

$$2^{52+e}|ad - bc| \leq 2^{2e}bd \leq ac. \quad (2)$$

Using (1) if $e \leq 0$ and (2) if $e \geq 0$, along with the bound on $\max(a, b, c, d)$, we have

$$2^{52+|e|}|ad - bc| \leq 67114657^2 < 2^{53}. \quad (3)$$

hence

$$2^{|e|}|ad - bc| < 2. \quad (4)$$

It follows that either $|ad - bc| = 0$ (in which case $a/b = c/d$ as required), or $|ad - bc| = 1$ and $e = 0$. In this case the inequality (1) gives $2^{52} \leq bd$, and the inequalities $2^e \leq a/b$ and $2^e \leq c/d$ give $b \leq a$ and $c \leq d$. At this point we look for a contradiction.

Swapping a/b and c/d if necessary, we can assume that $d \leq b$, so from $2^{52} \leq bd$ we have $2^{26} \leq b \leq a$. Now it's possible to do an exhaustive search over all pairs (a, b) of integers satisfying

$$2^{26} \leq b \leq a \leq 67114657.$$

There are exactly 16788115 such pairs (reducing to 10204542 after we discard those with $\gcd(a, b) \neq 1$), making this search computationally very feasible. For each pair (a, b) we can use the extended Euclidean algorithm to efficiently find all possible solutions in integers c and d to $|ad - bc| = 1$ with $0 < d \leq b$ (there are at most two), and check that a/b and c/d do not round to the same float, giving us our contradiction.

There remains the second case, where there is no integer e such that both a/b and c/d lie in $[2^e, 2^{e+1}]$. The only way for this to be possible is if there's a power of two separating a/b and c/d ; that is, without loss of generality (swapping a/b and c/d if necessary) there's an e satisfying

$$a/b < 2^e < c/d.$$

From the bounds on a , b , c and d , we must have $-26 \leq e \leq 26$. But now 2^e can be expressed as a fraction with numerator and denominator not exceeding 67114657, and since a/b and c/d round to the same float, 2^e (being squeezed between a/b and c/d) must round to that same float. So we can apply the case 1 proof to a/b and 2^e to deduce that $a/b = 2^e$, and again to 2^e and c/d to deduce that $2^e = c/d$, hence that $a/b = c/d$, as required. \square