fn get_min_k

Michael Shoemate

This proof resides in "contrib" because it has not completed the vetting process.

Proves soundness of the implementation of get_min_k in mod.rs at commit f5bb719 (outdated1).

1 Hoare Triple

Precondition

Compiler-Verified

- Generic T implements trait Float
- Type i32 implements the trait ExactIntCast<T.Bits>, where T.Bits is the type of the native bit representation of T.

User-Verified

None

Pseudocode

```
def get_min_k() -> i32:
return -i32.exact_int_cast(T.EXPONENT_BIAS) - i32.exact_int_cast(T.MANTISSA_BITS) + 1
```

Postcondition

Theorem 1.1. Return the k where 2^k is the smallest distance between adjacent non-equal values in T.

Proof. The floating-point exponent distinguishes between bands of floating-point numbers, where each band halves in width and halves the gap between adjacent numbers in the band. The smallest such band occurs when the exponent is smallest. This is the case where the unbiased exponent is all zeroes, and the biased exponent is simply the bias. This band of numbers is called the subnormals. The largest subnormal is almost 2^{1-b} , where b is the bias (b is 1023 for 64-bit floats).

This band of subnormal numbers is further sub-divided evenly into 2^m values by the m bits of the mantissa.

Together, the gap between the smallest adjacent non-zero values of type T is $2^{1-b} \cdot 2^{-m} = 2^{-b-m+1}$. This is implemented with associated constants for type T, where T.EXPONENT_BIAS is b and T.MANTISSA_BITS is m.

¹See new changes with git diff f5bb719..e2417bf5 rust/src/measurements/noise/nature/float/utilities/mod.rs