

Math 56 Notes for lec. 3 end, 1/4/14

FLOATING-POINT REPRESENTATION

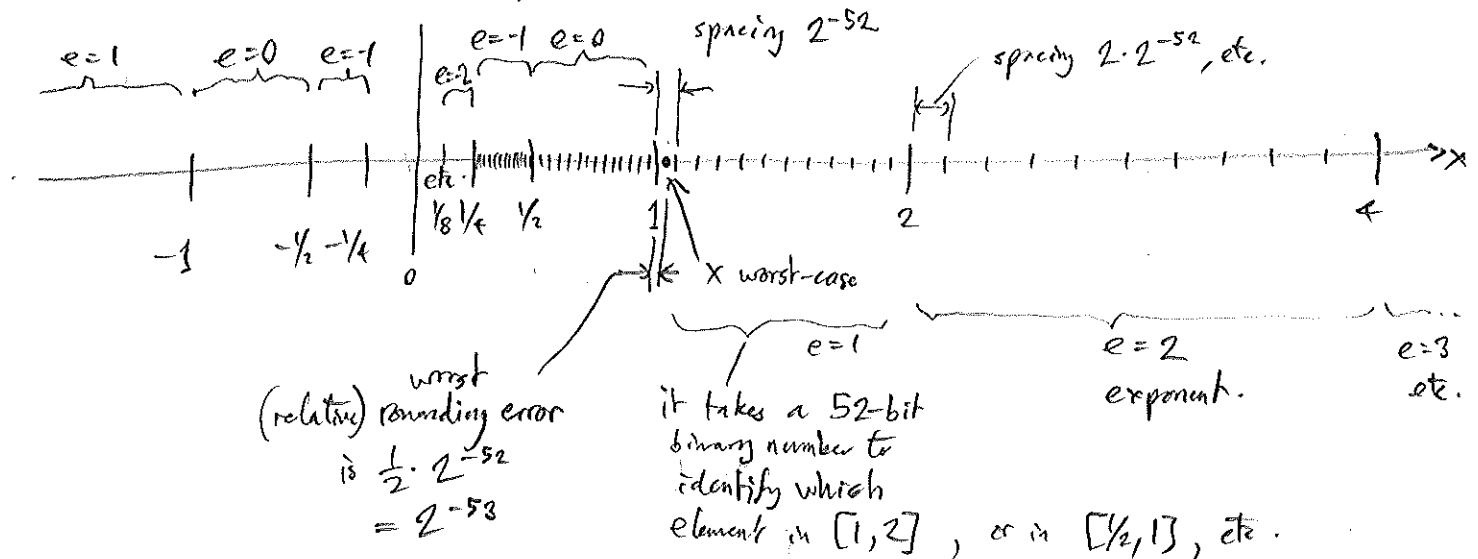
$fl(x)$ means the machine's representation of $x \in \mathbb{R}$, which is approximate.

It has small relative error in the sense that,

$$fl(x) = (1 + \epsilon)x \quad \text{for some } |\epsilon| \leq \epsilon_{\text{mach}}$$

For 'IEEE double precision' (the standard), $\epsilon_{\text{mach}} = 2^{-53} \approx 1.1 \times 10^{-16}$

How is this achieved? $fl(x) \in F$, a discrete set:



In general base β (here $\beta=2$), and precision t (here $t=52$)

So $x \in [1/\beta, 1]$ can be rep. by $fl(x) = \frac{m}{\beta^t}$ where $\beta^{t-1} \leq m \leq \beta^t$

To reach the rescaled copies of this, multiply by β^e where $e =$ "exponent", integer.
 (integer "mantissa")

Double-prec. devotes 11 bits to e , so $-1022 \leq e \leq 1023$

This means the largest & smallest sized numbers in F are $2^{1023} \approx 10^{308}$

The full set $F = \{0, \pm \frac{m}{\beta^t} \beta^e, \pm \infty, NaN\}$
 $\pm \infty$: not a number. $\pm NaN$: sign e in binary.
 2^{1023} & $2^{-1022} \approx 10^{-308}$

This is packed into 64 bits (8 bytes) as follows:
 1 bit (sign) | 11 bits (e) | 52 bits (m in binary)