

Quote:

"Floating point numbers are a lot like sandpiles:
every time you move one you lose a little sand
and pick up a little dirt."

- Koenigsmann + Plauger

-
- IEEE f.p. standard
 - Floating point arithmetic
 - Backward stability

IEEE floating pt standard:

$$x = 1 \cdot \underbrace{d_1 d_2 \dots d_t}_{\text{mantissa}} \beta^e \quad \begin{matrix} \nearrow \text{exponent} \\ \downarrow \text{base} \end{matrix}$$

β (base), typically 2 or 16 (2 for our purposes), $t = \beta - 1$.

Single precision: $\beta = 2$, $t = 23$, $-126 \leq e \leq 127$, total size: 32 bits
 \uparrow
 $2^7 - 1$

Double precision: $\beta = 2$, $t = 52$, $-1022 \leq e \leq 1023$, total size: 64 bits.
 \uparrow
 $2^{10} - 1$

\mathbb{F} : Floating pt numbers.

Machine ϵ ϵ_m is the smallest x s.t. $x+1 > 1$ in computer arithmetic.

Algorithm to find ϵ_m

```

e = 1
r = 2
while r > 1
    e = e/2
    r = 1 + e
end
 $\epsilon_m = 2e$ 
    
```

{converges in 53 iterations. Why?}

Notice: if $x \in \mathbb{R}$, $x = (1.d_1 d_2 \dots d_t d_{t+1} \dots) \times \beta^e$
 $\exists \hat{x} \in \mathbb{F}_2 \quad \hat{x} = (1.d_1 d_2 \dots d_t) \times \beta^e$
 $\tilde{x} \in \mathbb{F}_2 \quad \tilde{x} = (1.d_1 d_2 \dots d'_t) \times \beta^e$
 $|x - \hat{x}| < d_{t+1} \times \beta^e \rightarrow \text{dropping}$
 $|x - \tilde{x}| < \frac{1}{\beta} d_t \times \beta^e \rightarrow \text{rounding}$

Another source of error: operations on floating-pt numbers.

10f.5

exact operations: $+$, $-$, \times , \div

floating-pt ops: \oplus , \ominus , \otimes , \oslash

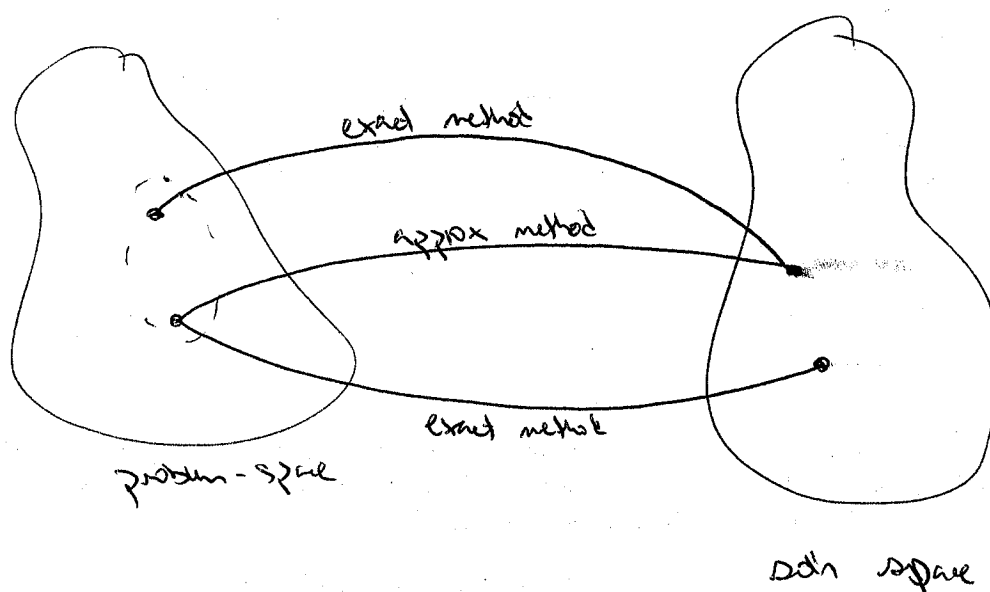
Fundamental axiom of F.P. arithmetic:

For all $x, y \in \mathbb{F}$, $\exists \delta$ with $|\delta| \leq \epsilon_m$ s.t.

$$x \otimes y = (x * y)(1 + \delta), \quad * = \{+, -, \times, \div\}.$$

\Rightarrow operation perturbations are on the order of ϵ_m .

Recall: Backwards stability



An algorithm is "backwards stable" if the approximate soln to a given problem is the exact soln of a nearby problem.

~~Next: analysis of complex multiply, back ops~~

Next: rounding error analysis

Stability of Algorithms

Stab-1

$f(y)$ \rightarrow exact soln of a problem f with data y

$\tilde{f}(y)$ \rightarrow inexact (numerical soln)

Algorithm is accurate if $\frac{\|f(y) - \tilde{f}(y)\|}{\|f(y)\|} \leq C \epsilon$ for ϵ small. (i.e. $\epsilon = \epsilon_m$)

Ex: Complex multiply:

ACT: $\underbrace{x \otimes y}_{\tilde{f}(x,y)} = \underbrace{(x \times y)}_{f(x,y)} (1 + \underbrace{2\sqrt{2}}_C \epsilon)$ where $|\epsilon| \leq \epsilon_m$

$$\frac{\|f(x,y) - \tilde{f}(x,y)\|}{\|f(x,y)\|} \leq \underbrace{2\sqrt{2}}_C \epsilon_m$$

Backwards stability:

Algorithm is backwards stable if $\tilde{f}(x) = f(\tilde{x})$ for some \tilde{x} satisfying

\tilde{x} exact soln to nearby problem

$$\frac{\|x - \tilde{x}\|}{\|x\|} \leq C \epsilon_m$$

$$x \otimes y = x \times (1 + 2\sqrt{2}\epsilon) y, \quad |\epsilon| \leq \epsilon_m$$

$$\text{Let } \tilde{x} = x, \quad \tilde{y} = (1 + 2\sqrt{2}\epsilon) y$$

$$\tilde{x} \times \tilde{y} = x \otimes y$$

$$\frac{\| \begin{pmatrix} x \\ y \end{pmatrix} - \begin{pmatrix} \tilde{x} \\ \tilde{y} \end{pmatrix} \|}{\| \begin{pmatrix} x \\ y \end{pmatrix} \|} = \frac{\| \begin{pmatrix} 0 \\ 2\sqrt{2}\epsilon y \end{pmatrix} \|}{\| \begin{pmatrix} x \\ y \end{pmatrix} \|} = \frac{2\sqrt{2}|\epsilon||y|}{\sqrt{x^2 + y^2}} \leq \underbrace{2\sqrt{2}}_C \epsilon_m$$