

Feb 19, 2021 Floating point
 Approx real numbers + arithmetic
 on computers

Error in (1) representation
 (2) arithmetic

FP is like binary scientific notation

sign (-) fraction (5.282) exponent (10)
 lead digit (5) base (10)

$$1.011 \cdot 2^{-11} = \left(1 + \frac{1}{4} + \frac{1}{8}\right) \cdot 2^{-(1+2)}$$

$$= 0.171875$$

Normalized FP assumes lead digit = 1

$$\pm 1.b_1 \dots b_s \cdot 2^p$$

significand

IEEE standard (64-bit)

1	1	52
---	---	----

sign s exp e fraction f $(-1)^s (1+f) \cdot 2^{e-1023}$

$2^{11} = 2048$

IEEE (32-bit)

1	8	23
---	---	----

$e - 127$

= value of s after iterations (\hat{s}_{n-1})

$$\hat{s}_0 = 0 = s_0 \quad \hat{s}_1 = x_1 = s_1$$

$$\hat{s}_k = f(\hat{s}_{k-1} + x_k) = (\hat{s}_{k-1} + x_k) \cdot (1 + \delta_k) \quad |\delta_k| \leq \epsilon$$

$$\hat{s}_2 = (x_1 + x_2)(1 + \delta_2)$$

$$\hat{s}_3 = (\hat{s}_2 + x_3)(1 + \delta_3) = x_1(1 + \delta_1 + \delta_2 + \delta_2\delta_3) + x_2(1 + \delta_1 + \delta_2 + \delta_2\delta_3) + x_3(1 + \delta_3)$$

$\overset{\sim}{\delta_2\delta_3} \sim 0$ $\overset{\sim}{\delta_2\delta_3} \sim 0$

$$\hat{s}_n = \underbrace{x_1(1 + \sum_{j=2}^n \delta_j)}_{\tilde{x}_1} + \underbrace{x_2(1 + \sum_{j=2}^n \delta_j)}_{\tilde{x}_2} + \underbrace{x_3(1 + \sum_{j=3}^n \delta_j)}_{\tilde{x}_3} + \dots + x_n(1 + \sum_{j=n}^n \delta_j)$$