

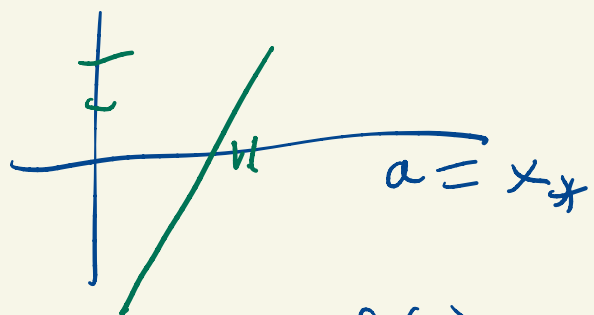
Floating Point

Math 426

University of Alaska Fairbanks

September 23, 2020

$$f(x) = f(a) + f'(a)(x-a) + \underbrace{O((x-a)^2)}_R$$



R

$$f(x) = 0 + f'(x_#)(x-x_#) + O((x-x_#)^2)$$

$$f'(x_#) = 0$$

$$\frac{f''(\xi)}{2} (x-x_#)^2$$

$$\underbrace{|f(x_k)|}$$

$$\approx \underbrace{|f'(x_#)| |x_k - x_#|}_{|e_k|}$$

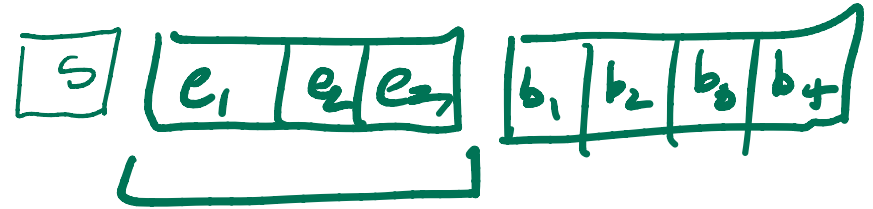
$$|e_k|$$

$$+ \underbrace{O(e_k^2)}$$

Last class

Fictional 8-bit floating point

1. 1 sign bit
2. 3 exponent bits (offset binary)
3. 4 mantissa bits (hidden bit representation)



$s_1 e_1 e_2 e_3 b_1 b_2 b_3 b_4$ represents

$$(-1)^{s_1} (1.b_1 b_2 b_3 b_4)_2 \times 2^E$$

where $E = (e_1 e_2 e_3)_2 - \Omega$; offset $\Omega = 2^2 - 1$

$$0 - (2^2 - 1) = 1 - 2^2$$

2^3 patterns

2^2 negative

2^2 positive

Last class

Fictional 8-bit floating point

1. 1 sign bit
2. 3 exponent bits (offset binary)
3. 4 mantissa bits (hidden bit representation)

$s_1 e_1 e_2 e_3 b_1 b_2 b_3 b_4$ represents

$$(-1)^{s_1} (1.b_1 b_2 b_3 b_4)_2 \times 2^E$$

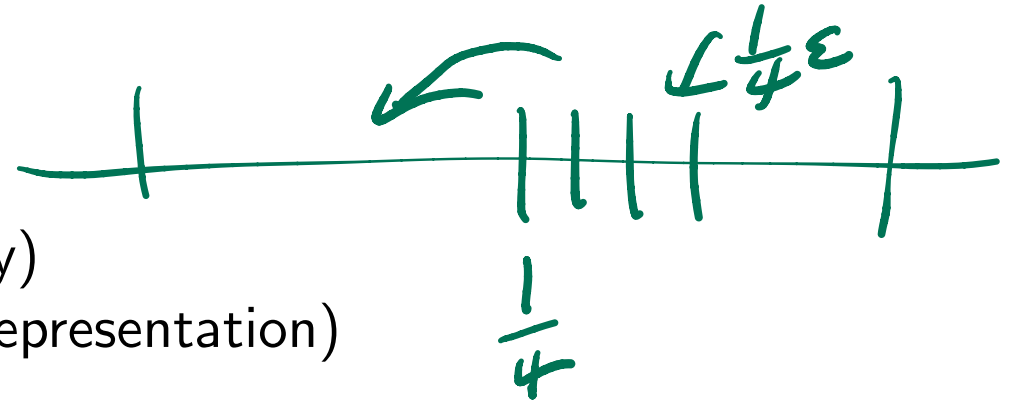
where $E = (e_1 e_2 e_3)_2 - \Omega$; offset $\Omega = 2^2 - 1$

Machine ϵ : 2^{-4}

Last class

Fictional 8-bit floating point

1. 1 sign bit
2. 3 exponent bits (offset binary)
3. 4 mantissa bits (hidden bit representation)



$s_1 e_1 e_2 e_3 b_1 b_2 b_3 b_4$ represents

$$(-1)^{s_1} (1.b_1 b_2 b_3 b_4)_2 \times 2^E$$

where $E = (e_1 e_2 e_3)_2 - \Omega$; offset $\Omega = 2^2 - 1$

Machine ϵ : 2^{-4}

$1, 1 + \epsilon$

except, if $e_1 e_2 e_3 = 000$:

$$(-1)^{s_1} (0.b_1 b_2 b_3 b_4)_2 \times 2^{(001)_2 - \Omega}$$

(two zeros; remainder are subnormal)

000
001
.
:
111

Last class

Fictional 8-bit floating point

1. 1 sign bit
2. 3 exponent bits (offset binary)
3. 4 mantissa bits (hidden bit representation)

$s_1 e_1 e_2 e_3 b_1 b_2 b_3 b_4$ represents

$$(-1)^{s_1} (1.b_1 b_2 b_3 b_4)_2 \times 2^E$$

where $E = (e_1 e_2 e_3)_2 - \Omega$; offset $\Omega = 2^2 - 1$

Machine ϵ : 2^{-4}

except, if $e_1 e_2 e_3 = 000$:

$$(-1)^{s_1} (0.b_1 b_2 b_3 b_4)_2 \times 2^{(001)_2 - \Omega}$$

(two zeros; remainder are subnormal) except, if $e_1 e_2 e_3 = 111$:

If $b_1 b_2 b_3 b_4 = 0000$, is $(-1)^s \infty$

Otherwise, is NaN.

IEEE 754

Single precision: 32 bits.

1. 1 sign bit
2. 8 exponent bits
3. 23 mantissa bits

Offset $\Omega = 2^7 - 1$.

$$(-1)^s (1.b_1 \dots b_{23}) 2^E$$

where $E = (e_1 \dots e_8)_2 - \Omega$.

$$2^8 / 2 = 2^7$$

$$2^7 \text{ neg.}$$

$$2^7 \text{ pos.}$$

$$-2^7 + 1$$

IEEE 754

Single precision: 32 bits.

1. 1 sign bit
2. 8 exponent bits
3. 23 mantissa bits

Offset $\Omega = 2^7 - 1$.

$$(-1)^s (1.b_1 \dots b_{23}) 2^E$$

where $E = (e_1 \dots e_8) - \Omega$. Machine ϵ ?

$$\begin{array}{l} (-1)^0 (1.0 \dots 0) 2^0 \\ (-1)^0 (1.0 \dots 01) 2^0 \end{array} \left. \vphantom{\begin{array}{l} (-1)^0 (1.0 \dots 0) 2^0 \\ (-1)^0 (1.0 \dots 01) 2^0 \end{array}} \right\} 2^{-23} = \epsilon$$

$\uparrow \uparrow \dots \uparrow$
1 2 23

IEEE 754

Single precision: 32 bits.

1. 1 sign bit
2. 8 exponent bits
3. 23 mantissa bits

Offset $\Omega = 2^7 - 1$.

$$(-1)^s (1.b_1 \cdots b_{23}) 2^E$$

where $E = (e_1 \cdots e_8) - \Omega$. Machine ϵ ? $\underline{2^{-23}} \approx \underline{1.1 \times 10^{-7}}$

IEEE 754

Single precision: 32 bits.

1. 1 sign bit
2. 8 exponent bits
3. 23 mantissa bits

$$1 - (2^7 - 1) = 2 - 2^7$$

Offset $\Omega = 2^7 - 1$.

$$(-1)^s \underbrace{(1.b_1 \dots b_{23})}_{2} 2^E$$

where $E = (e_1 \dots e_8) - \Omega$. Machine ϵ ? $2^{-23} \approx 1.1 \times 10^{-7}$

Smallest normal number?

positive

$$(-1)^0 (1.0 \dots 0) 2^E \rightarrow$$

IEEE 754

Single precision: 32 bits.

1. 1 sign bit
2. 8 exponent bits
3. 23 mantissa bits

Offset $\Omega = 2^7 - 1$.

$$(-1)^s (1.b_1 \cdots b_{23}) 2^E$$

where $E = (e_1 \cdots e_8) - \Omega$. Machine ϵ ? $2^{-23} \approx 1.1 \times 10^{-7}$

Smallest normal number?

$$2^{1-\Omega} = 2^{2-2^7} = 2^{-126} \approx 1.1 \times 10^{-38}$$

IEEE 754

Single precision: 32 bits.

1. 1 sign bit
2. 8 exponent bits
3. 23 mantissa bits

Offset $\Omega = 2^7 - 1$.

$$(-1)^s (1.b_1 \dots b_{23}) 2^E$$

where $E = (e_1 \dots e_8) - \Omega$. Machine ϵ ? $2^{-23} \approx 1.1 \times 10^{-7}$

Smallest normal number?

$$2^{1-\Omega} = 2^{2-2^7} = 2^{-126} \approx 1.1 \times 10^{-38}$$

Largest number?

$$(-1)^0 \overbrace{(1.1 \dots 1)}^{2-\epsilon} \cdot 2^E$$

11111110₂
100000000
(2⁸ - 2 - Ω)

IEEE 754

Single precision: 32 bits.

1. 1 sign bit
2. 8 exponent bits
3. 23 mantissa bits

Offset $\Omega = 2^7 - 1$.

$$(-1)^s (1.b_1 \cdots b_{23}) 2^E$$

where $E = (e_1 \cdots e_8) - \Omega$. Machine ϵ ? $2^{-23} \approx 1.1 \times 10^{-7}$

Smallest normal number?

$$2^{1-\Omega} = 2^{2-2^7} = 2^{-126} \approx 1.1 \times 10^{-38}$$

Largest number?

$$2^{(2^8-2)-\Omega} (2 - \epsilon) = 2^{2^8-2-2^7+1} = 2^{127} (2 - \epsilon) \approx 2^{128} \approx 3.4 \times 10^{38}$$

IEEE 754

~~Single~~ precision: 64 bits.

Double

1. 1 sign bit
2. 11 exponent bits
3. 52 mantissa bits

Offset $\Omega = 2^{10} - 1$.

$$(-1)^s (1.b_1 \cdots b_{52}) 2^E$$

where $E = (e_1 \cdots e_{11})_2 - \Omega$.

11

IEEE 754

Single precision: 64 bits.

1. 1 sign bit
2. 11 exponent bits
3. 52 mantissa bits

Offset $\Omega = 2^{10} - 1$.

$$(-1)^s (1.b_1 \cdots b_{52}) 2^E$$

where $E = (e_1 \cdots e_{11})_2 - \Omega$. Machine ϵ ?



2^{-52}

IEEE 754

Single precision: 64 bits.

1. 1 sign bit
2. 11 exponent bits
3. 52 mantissa bits

Offset $\Omega = 2^{10} - 1$.

$$(-1)^s (1.b_1 \cdots b_{52}) 2^E$$

where $E = (e_1 \cdots e_{11})_2 - \Omega$. Machine ϵ ? $2^{-52} \approx 2.2 \times 10^{-16}$

IEEE 754

Single precision: 64 bits.

1. 1 sign bit
2. 11 exponent bits
3. 52 mantissa bits

Offset $\Omega = 2^{10} - 1$.

$$(-1)^s (1.b_1 \cdots b_{52}) 2^E$$

where $E = (e_1 \cdots e_{11})_2 - \Omega$. Machine ϵ ? $2^{-52} \approx 2.2 \times 10^{-16}$

Smallest normal number?

IEEE 754

Single precision: 64 bits.

1. 1 sign bit
2. 11 exponent bits
3. 52 mantissa bits

Offset $\Omega = 2^{10} - 1$.

$$(-1)^s (1.b_1 \cdots b_{52}) 2^E$$

where $E = (e_1 \cdots e_{11})_2 - \Omega$. Machine ϵ ? $2^{-52} \approx 2.2 \times 10^{-16}$

Smallest normal number?

$$2^{1-\Omega} = 2^{2-2^{10}} = 2^{-1022} \approx 2.2 \times 10^{-308}$$

IEEE 754

Single precision: 64 bits.

1. 1 sign bit
2. 11 exponent bits
3. 52 mantissa bits

Offset $\Omega = 2^{10} - 1$.

$$(-1)^s (1.b_1 \cdots b_{52})_2 2^E$$

where $E = (e_1 \cdots e_{11})_2 - \Omega$. Machine ϵ ? $2^{-52} \approx 2.2 \times 10^{-16}$

Smallest normal number?

$$2^{1-\Omega} = 2^{2-2^{10}} = 2^{-1022} \approx 2.2 \times 10^{-308}$$

Largest number?

$$(2 - \epsilon) \cdot 2^{(2^{11} - 2) - \Omega}$$

IEEE 754

Single precision: 64 bits.

1. 1 sign bit
2. 11 exponent bits
3. 52 mantissa bits

Offset $\Omega = 2^{10} - 1$.

$$(-1)^s (1.b_1 \cdots b_{52}) 2^E$$

where $E = (e_1 \cdots e_{11})_2 - \Omega$. Machine ϵ ? $2^{-52} \approx 2.2 \times 10^{-16}$

Smallest normal number?

$$2^{1-\Omega} = 2^{2-2^{10}} = 2^{-1022} \approx 2.2 \times 10^{-308}$$

Largest number?

$$2^{(2^{11}-2)-\Omega} (2 - \epsilon) = 2^{1023} (2 - \epsilon) \approx 1.8 \times 10^{308}$$

Matlab demo

$$(1.0 \dots 01)_2 \times 2^{-1022} / 2$$

Machine epsilon: eps

showfloat: (s,e,m) bit patterns; I'll put it on the website.

$$\begin{array}{l} 0 \quad 0 \dots 01 \quad 0 - \quad \text{---} \quad 0 \\ 1.0 \text{---} \dots 0 \times 2^{-1022} \quad \Rightarrow \quad 2^{-1022} \\ \quad \quad \quad \uparrow \\ \quad \quad \quad 29 \\ 0.10 \text{---} \dots 0 \times 2^{-1022} \end{array}$$